

Don't Blame Disks for Every Storage Subsystem Failure

A Comprehensive Study of Storage Subsystem Failure Characteristics

Weihang Jiang, Chongfeng Hu, Yuanyuan Zhou and Arkady Kanevsky

April 2008 | version 1.0

Abstract

Building reliable storage systems becomes increasingly challenging as the complexity of modern storage systems continues to grow. Understanding storage failure characteristics is crucially important for designing and building a reliable storage system. While several recent studies have been conducted on understanding storage failures, almost all of them focus on the failure characteristics of one component – disks – and do not study other storage component failures.

This report analyzes the failure characteristics of storage subsystems. More specifically, we analyzed the NetApp AutoSupport logs collected from about 39,000 storage systems commercially deployed at various customer sites. The data set covers a period of 44 months and includes about 1,800,000 disks hosted in about 155,000 storage shelf enclosures. Our study reveals many interesting findings, providing useful guideline for designing reliable storage systems. Some of our major findings include: (1) Physical interconnects failures make up the largest part (27-68%) of storage subsystem failures, disk failures make up the second largest part (20-55%). Choices of disk types, shelf enclosure models and other components of storage subsystems contribute to the variability. (2) Each individual storage subsystem failure type and storage subsystem failure as a whole exhibit strong self-correlations. In addition, these failures exhibit “bursty” patterns. (3) Storage subsystems configured with redundant interconnects experience 30-40% lower failure rates than those with a single interconnect. (4) Spanning disks of a RAID group across multiple shelves provides a more resilient solution for storage subsystems than within a single shelf.

1 Introduction

1.1 Motivation

Reliability is a critically important issue for storage systems because storage failures can not only cause service downtime, but also lead to data loss. Building reliable storage systems becomes increasingly challenging as the complexity of modern storage systems grows into an unprecedented level. For example, the EMCTM Symmetrix DMX-4 can be configured with up to 2400 disks [8], the GoogleTM File System cluster is composed of 1000 storage nodes [9], and the NetApp[®] FAS6000 series can support more than 1000 disks per node, with up to 24 nodes in a system [13].

To make things even worse, disks are not the only component in storage systems. To connect and access disks, modern storage systems also contain many other components, including shelf enclosures, cables and host adapters, and complex software protocol stacks. Failures in these components can lead to downtime and/or data loss of the storage system. Hence, in complex storage systems, component failures are very common and critical to storage system reliability.

To design and build a reliable storage system, it is crucially important to understand the storage failure characteristics. First, accurate estimation of storage failure rate can help system designers decide how many resources should be used to tolerate failures and to meet certain service-level agreement (SLA) metrics (*e.g.*, data availability). Second, knowledge about factors that greatly impact the storage system reliability can guide designers to select more reliable components or build redundancy into unreliable components. Third,

understanding the statistical properties such as failure distribution over time of modern storage systems is necessary to build right testbed and fault injection models to evaluate existing resiliency mechanisms and to develop better fault-tolerant mechanisms.

While several recent studies have been conducted on understanding storage failures, almost all of them focused on the failure characteristics of one storage component—disks. For example, disk vendors have studied the disk failure characteristics through running accelerated life tests and collecting statistics from their return unit databases [4, 22]. Based on such tests, they calculate the *mean-time-to-failure (MTTF)* and record it in a disk specification. For most of the disks, the specified MTTF is typically more than one million hours, equivalent to a lower than 1% *annualized failure rate (AFR)*. But such low AFR is usually not what has been experienced by users. Motivated from this observation, recently some researchers have studied *disk failures* from a user’s perspective by analyzing disk replacement logs collected in the field [15, 17]. Interestingly, they found disks are replaced much more frequently (2-4 times) than vendor-specified AFRs. But as this study indicates, there are other storage subsystem failures besides disk failures that are treated as disk faults and lead to unnecessary disk replacements. Additionally, some researchers analyzed the characteristics of disk sector errors, which can potentially lead to complete disk failures [2], and they found that sector errors exhibit strong temporal locality (*i.e.*, bursty patterns).

While previous works provide very good understanding of disk failures and an inspiring starting point, it is not enough since, besides disks, there are many other components that may contribute to storage failures. Without a good understanding of these components’ failure rates, failure distributions, and other characteristics, as well as impacts of these component failures on the storage system, it can make our estimation of the storage failure rate/distribution inaccurate. For example, as we will show in our study from real-world field data, having a lower disk failure rate does not necessarily mean that the corresponding storage system is more reliable—because some other components may not be as reliable.

More importantly, if we only focus on disk failures and ignore other component failures, we may fail to build a highly reliable storage system. For example, RAID is usually the primary resiliency mechanism built in to most modern storage systems (various forms of checksumming are considered as part of RAID). As RAID is mainly designed to tolerate disk failures, it is insufficient to handle other component failures such as failures in shelf enclosures, interconnects, and software protocol layers.

While we are interested in failures of a whole storage system, this study is concentrated on the core part of it — *the storage subsystem*, which contains disks and all components providing connectivity and usage of disks to the entire storage system.

We conducted a study using real-world field data from NetAppTM AutoSupport Database, to answer the following questions:

- How much do disk failures contribute to storage subsystem failures? What are other major factors that can lead to storage subsystem failures?
- What are the failure rates of other types of storage subsystem components such as physical interconnects and protocol stacks? What are the failure characteristics such as failure distribution and failure correlation for these components?
- Typically, some resiliency mechanisms such as RAID and redundancy mechanisms such as multipathing are used in practice to achieve high reliability and availability [5, 9]. How effective are these mechanisms in handling storage subsystem failures?

Data from the same AutoSupport Database was first analyzed in [2] on latent sector errors and was further analyzed in [3] on data corruptions.

There are other redundancy and resiliency mechanisms in storage system layers higher than the storage subsystem and RAID-based resiliency mechanism studied in this report. These mechanisms handle many storage subsystem failures. Studying impacts of these resiliency and redundancy mechanisms on storage

failures, including storage subsystem failures, is part of the future work.

1.2 Our Contributions

This report analyzes the failure characteristics of storage subsystems, including disks and other system components, based on a significant amount of field data collected from customers. Specifically, we analyzed the NetApp AutoSupport logs collected from about 39,000 storage systems commercially deployed at various customer sites. The data set covers a period of 44 months and includes about 1,800,000 disks hosted in about 155,000 storage shelf enclosures. Furthermore, our data covers a wide range of storage system classes, including *nearline (secondary)*, *low-end*, *mid-range*, and *high-end* systems.

This report studies failure characteristics from several angles. First, we classify storage subsystem failures into four failure types based on their symptoms and root causes and examine the relative frequency of each failure type. Second, we study the effect of several factors on storage subsystem reliability. These factors include disk models, shelf enclosure models, network redundancy mechanisms, and disk positioning in the shelf. Finally, we analyze the statistical properties of storage subsystem failures, including the correlation between failures and their time distribution.

Our study reveals many interesting findings, providing useful guideline for designing reliable storage systems. Following is a summary of our major findings and the corresponding implications:

- In addition to disk failures that contribute to 20-55% of storage subsystem failures, other components such as physical interconnects (including shelf enclosures) and protocol stacks also account for significant percentages (27-68% and 5-10%, respectively) of failures. Due to these component failures, even though storage systems of certain types (*e.g.*, low-end primary systems) use more reliable disks than some other types (*e.g.*, nearline systems), their storage subsystems exhibit higher failure rates. These results indicate that, to build highly reliable and available storage systems, only using resiliency mechanisms targeting disk failures (*e.g.*, RAID) is not enough. We also need to build resiliency mechanisms such as redundant physical interconnects and self-checking protocol stacks to tolerate failures in these storage components.
- Each individual storage subsystem failure type and storage subsystem failure as a whole exhibit strong correlations, (*i.e.* after one failure, the probability of additional failures of the same type is higher). In addition, failures also exhibit bursty patterns in time distribution, (*i.e.* multiple failures of the same type tend to happen relatively close together). These results motivate a revisiting of current resiliency mechanisms such as RAID that assume independent failures. These results also motivate development of better resiliency mechanisms that can tolerate multiple correlated failures and bursty failure behaviors.
- Storage subsystems configured with two independent interconnects experienced much (30-40%) lower AFRs than those with a single interconnect. This result indicates the importance of interconnect redundancy in the design of reliable storage systems.
- RAID groups built with disks spanning multiple shelf enclosures show much less bursty failure patterns than those built with disks from the same shelf enclosure. This indicates that the former is a more resilient solution for large storage systems.
- There is not prominent evidence that disk positioning in the shelf could affect its AFR.

The rest of the report is organized as follows. Section 2 provides the background and describes our methodology. Section 3 presents the contribution of disk failures to storage subsystem failures and frequency of other types of storage subsystem failures. Section 4 quantitatively analyzes the effects of several factors on storage subsystem reliability, while Section 5 analyzes the statistical properties of storage subsystem failures. Section 6 discusses the related work, and Section 7 concludes the report and provides directions for future work.

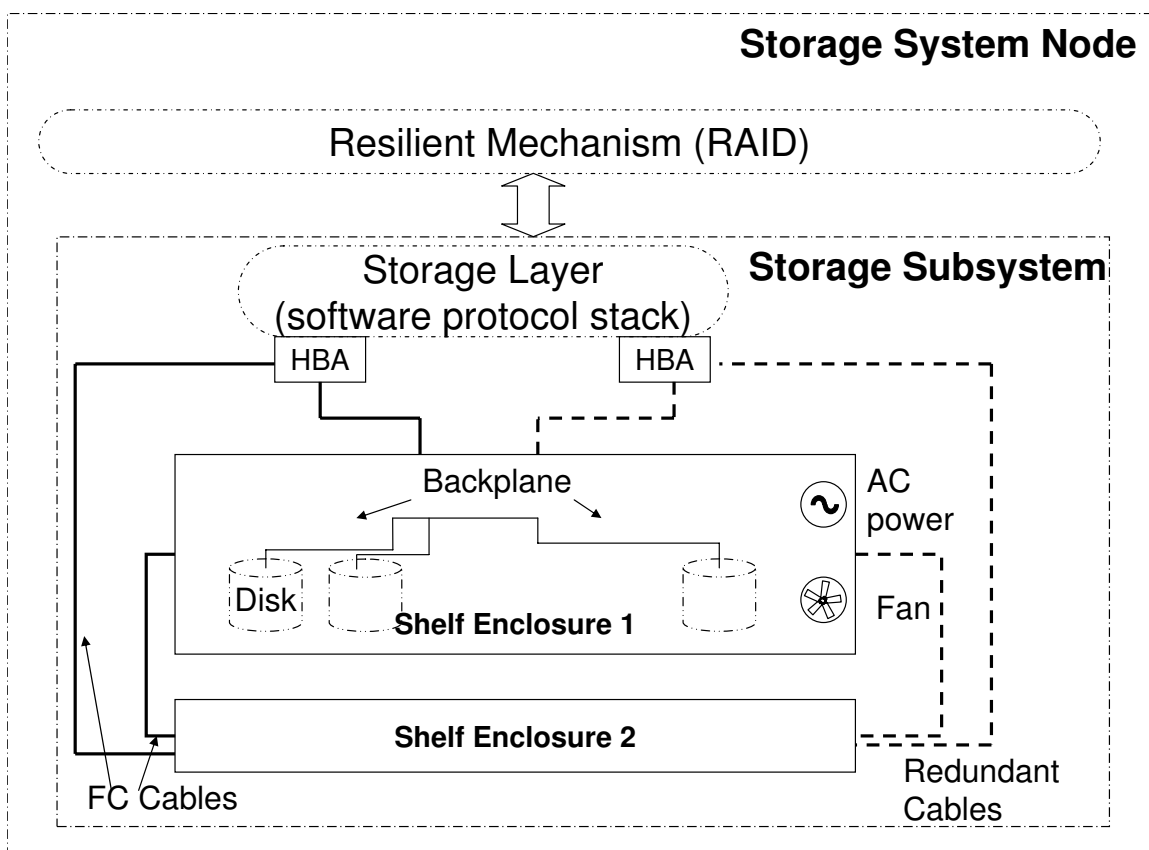


Figure 1. Storage system node architecture.

2 Background

In this section, we detail the typical architecture of storage systems we study in NetApp, the definition and terminology used in this report, and the source of the data studied in this report.

2.1 Storage System Architecture

Figure 1 shows the typical architecture of a NetApp storage system node. A NetApp storage system can be composed of several storage system nodes.

From the customers' perspective, a storage system is a virtual device that is attached to customers' systems and provides customers with the desired storage capacity with high reliability, good performance, and flexible management.

Looking from inside, a storage system node is composed of storage subsystems, resiliency mechanisms, storage head/controller, and other higher-level system layers. The storage subsystem is the core part of a storage system node and provides connectivity and usage of disks to the entire storage system node. It contains various components, including disks, shelf enclosures, cables and host adapters, and complex software protocol stacks. Shelf enclosures provide power supply, cooling service and prewired backplane for the disks mounted in them. Cables initiated from host adapters connect one or multiple shelf enclosures to

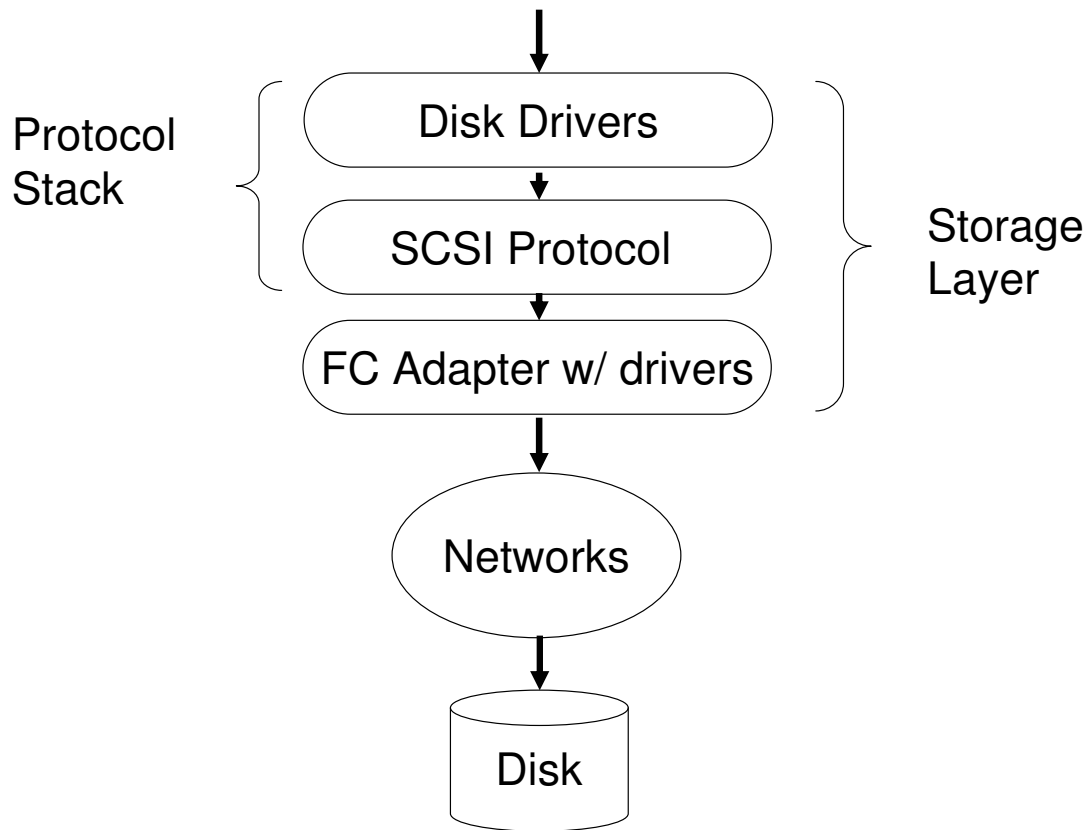


Figure 2. I/O request path in storage subsystem.

the network. Each shelf enclosure can be optionally connected to a secondary network for redundancy. In Section 4.3 we will show the impact of this redundancy mechanism on failures of the storage subsystem.

Usually, on top of the storage subsystem, resiliency mechanisms, such as RAID, are used to tolerate failures in storage subsystems.

2.2 Terminology

We use the followings terms in this report.

- **Disk family:** A particular disk product. The same product may be offered in different capacities. For example, “Seagate Cheetah 10k.7” is a disk family.
- **Disk model:** The combination of a disk family and a particular disk capacity. For example, “Seagate Cheetah 10k.7 300 GB” is a disk model. For disk family and disk model, we use the same naming convention as in [2, 3].
- **Failure types:** Refers to the four types of storage subsystem failures: disk failure, physical interconnect failure, protocol failure, and performance failure.
- **Shelf enclosure model:** A particular shelf enclosure product. All shelf enclosure models studied in this report can host at most 14 disks.

- **Storage subsystem failure:** Refers to failures that prevent the storage subsystem from providing storage service to the whole storage system node. However, not all storage subsystem failures are experienced by customers, since some of the failures can be handled by resiliency mechanisms on top of storage subsystems (*e.g.* RAID) and other mechanisms at higher layers.
- **Storage system class:** Refers to the capability and usage of storage systems. There are four storage system classes studied in this report: nearline systems (mainly used as secondary storage), low-end, mid-range, and high-end (mainly used as primary storage).
- Other terms in the report are used as defined by SNIA [20].

2.3 Definition and Classification of Storage Subsystem Failures

Figure 2 shows the steps and components that are involved in fulfilling an I/O request in a storage subsystem. As shown in Figure 2, for the storage layer to fulfill an I/O request, the I/O request will first be processed and transformed by protocols and then delivered to disks through networks initiated by host adapters. *Storage subsystem failures* are the failures that break the I/O request path, and can be caused by hardware failures, software bugs, and protocol incompatibilities along the path.

To better understand storage subsystem failures, we categorize them into four types along the I/O request path:

- **Disk Failure:** This type of failure is triggered by failure mechanisms of disks. Imperfect media, media scratches caused by loose particles, rotational vibration, and many other factors internal to a disk can lead to this type of failures. Sometimes, the storage layer proactively fails disks based on statistics collected by on-disk health monitoring mechanisms (*e.g.*, a disk has experienced too many sector errors [1]). These incidences are also counted as *disk failures*.
- **Physical Interconnect Failure:** This type of failure is triggered by errors in the networks connecting disks and storage heads. It can be caused by host adapter failures, broken cables, shelf enclosure power outage, shelf backplanes errors, and/or errors in shelf FC drivers. When *physical interconnect failures* happen, affected disks appear to be missing from the system.
- **Protocol Failure:** This type of failure is caused by incompatibility between protocols in disk drivers or shelf enclosures and storage heads and software bugs in the disk drivers. When this type of failure happens, disks are visible to the storage layer but I/O requests are not correctly responded by disks.
- **Performance Failure:** This type of failure happens when the storage layer detects that a disk cannot serve I/O requests in a timely manner while none of previous three types of failures are detected. It is mainly caused by partial failures, such as unstable connectivity or when disks are heavily loaded with disk-level recovery (*e.g.*, broken sector remapping).

The occurrences of these four types of failures are recorded in AutoSupport logs collected by NetApp.

2.4 Data Sources

Table 1 provides an overview of the data used in this study. NetApp AutoSupport logs from about 39,000 commercially deployed storage systems in four system classes are used for the results presented in this report. There are totally about 1,800,000 disks mounted in 155,000 shelf enclosures. The disks are a combination of SATA and FC disks. The population of disks contains at least 9 disk families and 15 disk models. The NetApp AutoSupport logs used for this study were collected between January 2004 and August 2007.

Below we describe each storage system class.

Nearline systems are deployed as cost-efficient archival or secondary storage systems. Less expensive SATA disks are used in nearline systems. In nearline systems, one storage subsystem on average contains about 7

System Classes	Duration	# Systems	# Shelves	Multipathing	# Disks	Disk Types	# RAID Groups	RAID Types	# Failure Types	# Failure Events
Nearline	1/04 - 8/07	4,927	33,681	single path	520,776	SATA	67,227	RAID4 RAID6	Disk Failure Phy. Inter. Failure Protocol Failure Performance Failure	10,105 4,888 1,819 1,080
Low-end	1/04 - 8/07	22,031	37,260	single-path	264,983	FC	44,252	RAID4 RAID6	Disk Failure Phy. Inter. Failure Protocol Failure Performance Failure	3,230 4,338 1,021 1,235
Mid-range	1/04 - 8/07	7,154	52,621	single-path dual-path	578,980	FC	77,831	RAID4 RAID6	Disk Failure Phy. Inter. Failure Protocol Failure Performance Failure	8,989 7,949 2,298 2,060
High-end	1/04 - 8/07	5,003	33,428	single-path dual-path	454,684	FC	49,555	RAID4 RAID6	Disk Failure Phy. Inter. Failure Protocol Failure Performance Failure	8,240 7,395 1,576 153

Table 1. Overview of studied storage systems. Note that the “# Disks” given in the table is the number of disks that have ever been installed in the system during the 44 months. For some systems, disks have been replaced during the period, and we account for that in our analysis by calculating the life time of each individual disk. The “# Failure Events” given in the table are the numbers of the four types of storage subsystem failures (disk failure, physical interconnect failure, protocol failure, and performance failure) that happened during the period.

shelf enclosures and 98 disks. Both RAID4 and RAID6 are supported as resiliency mechanisms in nearline systems.

Primary storage systems, including low, mid, and high-end systems, are mainly used in mission- or business-critical environments and primarily use FC disks. Low-end storage systems have embedded storage heads with shelf enclosures, but external shelf enclosures can be added. Mid-range and high-end systems use external shelves and are usually configured with more shelf enclosures and disks than low-end systems. Each mid-range system has about 7 shelf enclosures and 80 disks (not every shelf is fully utilized and configured with 14 disks), and high-end systems are in similar scale. Going from low to high-end systems, more reliable components and more redundancy mechanisms are used. For example, both mid-range and high-end systems support dual paths for redundant connectivity.

2.5 AutoSupport Logs and Analysis

The storage systems studied in this report have a low-overhead logging mechanism that automatically records informational and error events on each layer (software and hardware) and each subsystem during operation. Several recent works such as [2, 3] also studied the same set of NetApp AutoSupport logs from different aspects.

Figure 3 shows a log example that reports a physical interconnect failure. As can be seen in the figure, when a failure happens, multiple events are generated as the failure propagates from lower layers to higher layers (Fibre Channel to SCSI to RAID). By keeping track of events generated by lower layers, higher layers can identify the cause of events and tag the events with corresponding failure types. In this example, the RAID layer, which is right above the storage subsystem, generates a disk missing event, indicating a physical interconnect failure. In this report, we look at four types of events generated by the RAID layer, corresponding to four types of storage subsystem failures.

Besides the events shown in the example, there are many other events recorded in the NetApp AutoSupport logs. For example, standard error reports from the SCSI protocol layer tell us what failure mechanisms happen inside disks [19]. Disk medium error messages from disk drivers provide information about broken sectors [2]. Similarly, messages from FC protocol and FC host adapter drivers report errors that occur in FC networks and FC adapters.

It is important to notice that not all failures propagate to the RAID layer, as some failures are recovered or tolerated by storage subsystems. For example, an interconnect failure can be recovered through retries at SCSI layer or be tolerated through multipathing. Therefore, storage failures characterized as storage subsystem failure as a whole are those errors exposed by storage subsystems to the rest of the system.

- Sun Jul 23 05:43:36 PDT [fci.device.timeout:error]: Adapter 8 encountered a device timeout on device 8.24
- Sun Jul 23 05:43:50 PDT [fci.adapter.reset:info]: Resetting Fibre Channel adapter 8.
- Sun Jul 23 05:43:50 PDT [scsi.cmd.abortedByHost:error]: Device 8.24: Command aborted by host adapter:
- Sun Jul 23 05:44:12 PDT [scsi.cmd.selectionTimeout:error]: Device 8.24: Adapter/target error: Targeted device did not respond to requested I/O. I/O will be retried.
- Sun Jul 23 05:44:22 PDT [scsi.cmd.noMorePaths:error]: Device 8.24: No more paths to device. All retries have failed.
- Sun Jul 23 05:46:22 PDT [raid.config.filesystem.disk.missing:info]: File system Disk 8.24 S/N [3EL03PAV00007111LR8W] is missing.

Figure 3. Example of a piece of log reporting a physical interconnect failure.

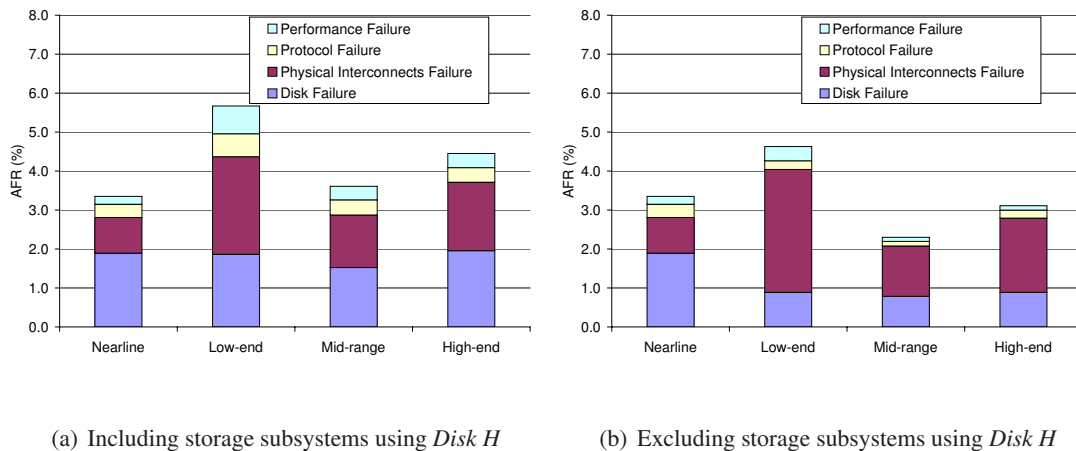


Figure 4. AFR for storage subsystems in four system classes and the breakdown based on failure types.

As Figure 3 shows, each event is tagged with the time-stamps when the failure is detected and with the ID of the disk affected by the failure. Since all the storage systems studied in this report periodically send data verification requests to all disks as a proactive method to detect failures, the lag between the occurrence and the detection of the failure is usually shorter than an hour.

System information is also copied with snapshots and recorded in NetApp AutoSupport logs on a weekly basis. This information is particularly important for understanding storage subsystem reliability since it provides the insight into the system parameters of storage subsystems. More specifically, NetApp AutoSupport logs contain the information about hardware components used in storage subsystems, such as disk models and shelf enclosure models, and they also contain the information about the layout of disks, such as which disks are mounted in the same shelf enclosures, and which disks are in the same RAID group. This information is used for analyzing statistical properties of storage subsystem failures in Section 5.

3 Frequency of Storage Subsystem Failures

As we categorize storage subsystem failures into four failure types based on their root causes, a natural question is therefore what the relative frequency of each failure type is. To answer this question, we study the NetApp AutoSupport logs collected from 39,000 storage systems.

Figure 4(a) presents the breakdown of AFR for storage subsystems based on failure types, for all four system classes studied in this report. Since one problematic disk family, denoted as *Disk H*, has already been reported in [2], for Figure 4(b) we exclude data from storage subsystems using *Disk H*, so that we can analyze the trend without being skewed by one problematic disk family. The discussion on *Disk H* is presented in Section 4.1.

Finding (1): *Physical interconnects failures* make up the largest part (27-68%) of storage subsystem failures, *disk failures* make up the second largest part (20-55%). *Protocol failures* and *performance failures* both make up noticeable fractions.

Implications: *Disk failures* are not always a dominant factor of storage subsystem failures, and a reliability study for storage subsystems cannot only focus on *disk failures*. Resilient mechanisms should target all failure types.

As Figure 4(b) shows, across all system classes, *disk failures* do not always dominate storage subsystem

failures. For example, in low-end storage systems, the AFR for storage subsystems is about 4.6%, while the AFR for *disks* is only 0.9%, about 20% of overall AFR. On the other hand, *physical interconnect failures* account for a significant fraction of storage subsystem failures, ranging from 27% to 68%. The other two failure types, *protocol failures* and *performance failures*, contribute to 5-10% and 4-8% of storage subsystem failures, respectively.

Finding (2): For *disks*, nearline storage systems show higher (1.9%) AFR than any primary storage systems (0.9%). But in spite of more reliable disks, low-end primary storage systems show higher (4.6%) AFR than nearline storage systems (3.4%).

Implications: *Disk failure* rate is not indicative of the storage subsystem failure rate.

Figure 4(b) also shows that nearline systems, which mostly use SATA disks, experience about 1.9% AFR for *disks*, while for low-end, mid-range, and high-end systems, which mostly use FC disks, the AFR for *disks* is under 0.9%. This observation is consistent with the common belief that enterprise disks (FC) are more reliable than nearline disks (SATA).

However, the AFR for storage subsystems does not follow the same trend. Storage subsystem AFR of nearline systems is about 3.4%, lower than that of low-end systems (4.6%). This indicates that other factors, such as shelf enclosure model and network configurations, strongly affect storage subsystem reliability. The impacts of these factors are examined in the next section.

Another interesting observation that can be seen in Figure 4(b) is that for FC drives, the *disk failure rate* is consistently below 1%, as published by disk drive manufacturers, while some previous works claim that the AFR for disks is much higher [15, 17]. We believe that the main reason for the discrepancy is that these studies look at disk failures from different angles. Our study is from a system's perspective, as we extract disk failure events from system logs, similar to disk drive manufacturers' studies. On the other hand [15, 17], look at disk failures from a user's perspective. Since their studies are based on disk replacement logs, they cannot identify the reasons for disk replacement. As system administrators often replace disks when they observe unavailability of disks, the disk replacement rates reported in these studies are actually close to the *storage subsystem failure rate* of this report.

4 Impact of System Parameters on Storage Subsystem Failures

As we have seen above, storage subsystems of different system classes show different AFRs. While these storage subsystems are architecturally similar, the characteristics of their components, like disks and shelves, and their redundancy mechanisms, like multipathing, differ. We now explore the impact of these factors on storage subsystem failures.

4.1 Disk Model

The disk is the key component of a storage subsystem; therefore it is important to understand how disk models affect storage subsystem failures. To understand the impact of the disk model, we study data collected from nearline, low-end, mid-range, and high-end systems.

Figure 5 shows the AFRs for storage subsystems from 4 system classes configured with 3 shelf enclosure models, 6 combinations in total (not every shelf enclosure model works with all system classes). Since we find that the enclosure model also has a strong impact on storage subsystem failures, we group data based on system class, shelf enclosure model, and disk model so that we can separately study the effects of these factors. In this section, we mainly focus on disk model; shelf enclosure model will be discussed in Section 4.2.

There are a total of 20 disk models used in these systems, and each disk model is denoted as *family-type*, with the same convention as in [2]. For anonymization purpose, a single letter is used to represent a disk

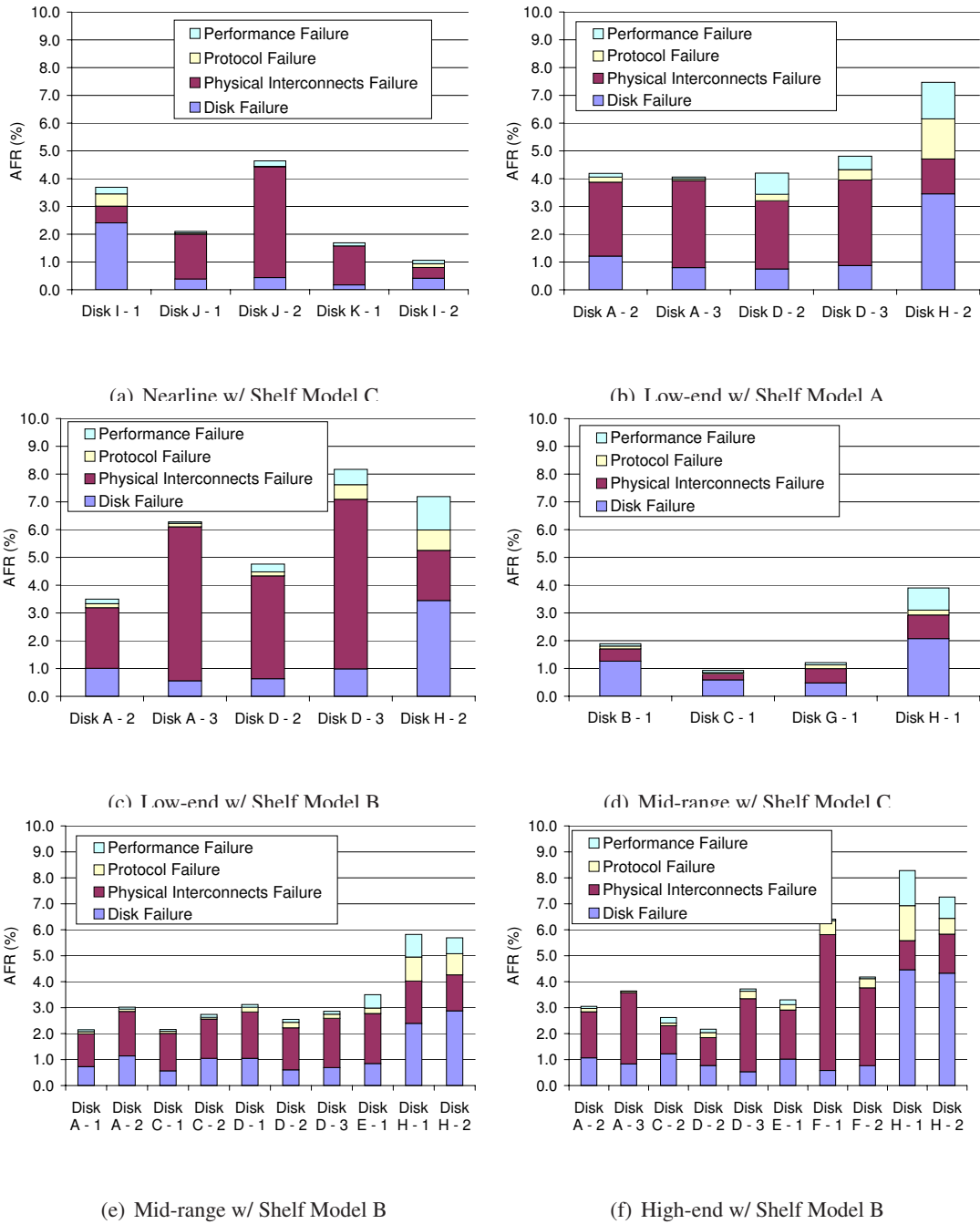


Figure 5. AFR for storage subsystems by disk models.

family (e.g., Seagate Cheetah 10k.7), and type is a single number indicating the disk's capacity. The relative capacity within a family is ordered by the number. For example, *Disk A-2* is larger than *A-1* and *B-2* is larger than *B-1*.

Finding (3): Storage subsystems using disks from a problematic disk family, show much higher (2 times) AFR than other storage subsystems.

Implications: Disk model is a critical factor to consider for designing reliable storage subsystems.

We can see from Figure 5 (a)-(f) that for most storage subsystems, AFR is about 2% - 4%. However, storage subsystems using *Disk H-1* and *Disk H-2* show 3.9%-8.3% AFR, higher than the average AFR by a factor of two.

We know that *Disk H-1* and *Disk H-2* are problematic. It is interesting to observe that not only disk failures but also protocol failures and performance failures are negatively affected by the problematic disks. The possible reason is that as disks experience failures, corner-case bugs in the protocol stacks are more likely to be triggered, leading to more occurrences of protocol failures. At the same time, some I/O requests cannot be served in time, causing more performance failures.

Finding (4): Storage subsystems using disks from the same disk models exhibit similar *disk failure* rates across different system environments (different system class or shelf enclosure models), but they show very different *storage subsystem failure* rates.

Implications: Factors other than disk models also heavily affect storage subsystem failures, while they are not revealed by *disk failures*.

As Figure 5 shows, some disk models are used by storage subsystems of multiple system classes, together with various shelf enclosure models. For example, *Disk A-2* and *Disk D-2* are used in low-end systems with different shelf models and by mid-range and high-end systems with the same shelf model.

As we can see from Figure 5, for the storage subsystems using the same disk models, *disk failure* rates do not change much. For example, *disk AFR* of *Disk D-2* varies from 0.6% to 0.77% with a standard deviation of 8%. For all storage subsystems sharing the same disk models, the average standard deviation of *disk AFR* is less than 11%.

On the other hand, the *storage subsystem AFR* exhibits strong variation. For example, AFR for *storage subsystems* using *Disk D-2* varies from 2.2% to 4.9%, with a standard deviation of 127%. For all such storage subsystems, the average standard deviation of *storage subsystem AFR* is as high as 98%. This observation indicates that *storage subsystem AFR* is strongly affected by factors other than *disk model*, while these factors do not affect *disk failures* much.

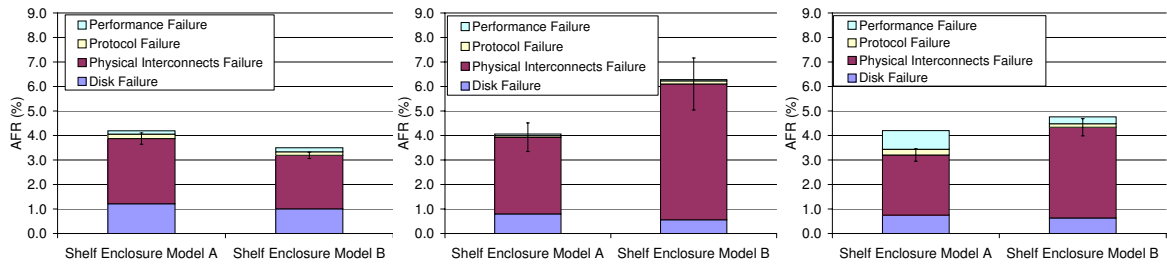
Finding (5): The AFR for *disks* and *storage subsystems* does not increase with disk size.

Implications: As disk capacity rapidly increases, storage subsystems will not necessarily experience more *disk failures* or *storage subsystem failures*.

We do not observe increasing *disk failure* rate or *storage subsystem failure* rate with increasing disk capacity. For example, as Figure 5 (e) shows, storage subsystems using *Disk D-2* show lower *disk* and *storage subsystem AFR* than those using *Disk D-1*.

4.2 Shelf Enclosure Model

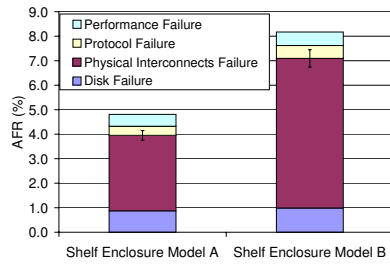
Shelf enclosures contain power supplies, cooling devices, and prewired backplanes that carry power and I/O bus signals to the disks mounted in them. Different shelf enclosure models are different in design and have different mechanisms for providing these services; therefore, it is interesting to see how shelf enclosure model affects storage subsystem failures.



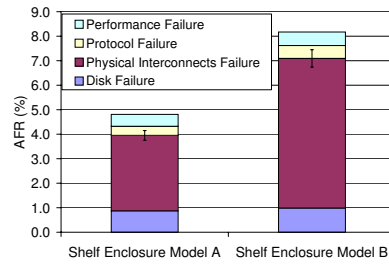
(a) Low-end Disk-A2

(b) Low-end Disk-A3

(c) Low-end Disk-D2



(d) Low-end Disk-D3



(e) Mid-range Disk-C1

Figure 6. AFR for storage subsystems by shelf enclosure models using the same disk models (a subset of data from Figure 5). The error bars show 99.5%+ confidence intervals for physical interconnect failures.

In order to study the impact of the shelf enclosure model, we look at the data collected from low-end storage systems, since low-end systems use the same disk models with different shelf enclosure models, so that we can study the effect of shelf enclosure models without inference from disk models.

Finding (6): The shelf enclosure model has a strong impact on storage subsystem failures, and different shelf enclosure models work better with different disk models.

Implications: To build a reliable storage subsystem, hardware components other than disks (*e.g.*, shelf enclosure) should also be carefully selected. And due to component interoperability issues, there might be a different “best choice” for one component depending on the choice of other components.

Figure 6 (a)-(e) shows AFR for storage subsystems when configured with different shelf enclosure models but the same disk models. As expected, shelf enclosure model primarily impacts *physical interconnect failures*, with little impact on other failure types, different from disk model, which impacts all failure types.

To confirm this observation, we tested the statistical significance using a T-test [16]. As Figure 6 (a) shows, the *physical interconnect failures* with different shelf enclosure models are quite different ($2.66 \pm 0.23\%$ versus $2.18 \pm 0.13\%$). A T-test shows that this is significant at the 99.5% confidence interval, indicating that the hypothesis that *physical interconnect failures* are impacted by shelf enclosure models is very strongly supported by the data. Figure 6(b)-(e) shows similar observations with significance at 99.5%, 99.9%, 99.9%, and 99.9% confidence.

It is also interesting to observe that for different disk models, different shelf enclosure models work better. For example, for *Disk-A2*, storage subsystems using *Shelf Enclosure B* show better reliability than those using *Shelf Enclosure A*, while for *Disk-A3*, *Disk-D2*, and *Disk-D3*, *Shelf Enclosure A* is more reliable. Such observations might be due to component interoperability issues between disks and shelf enclosures. This indicates that we might not be able to make the best decision on selecting the most reliable hardware components without evaluating the components from a system perspective and taking the effect of interoperability into account.

4.3 Network Redundancy Mechanism

As we have seen, *physical interconnect failures* contribute to a significant fraction (27-68%) of storage subsystem failures. Since *physical interconnect failures* are mainly caused by network connectivity issues in storage subsystems, it is important to understand the impact of network redundancy mechanisms on storage subsystem failures.

For the mid-range and high-end systems studied in this report, FC drivers support a network redundancy mechanism, commonly called *active/passive multipathing*. This network redundancy mechanism connects shelves to two independent FC networks, and redirects I/O requests through the redundant FC network when one FC network experiences network component failures (*e.g.*, broken cables).

To study the effect of this network redundancy mechanism, we look at the data collected from mid-range and high-end storage systems, and group them based on whether the network redundancy mechanism is turned on. As we observed from our data set, about 1/3 of storage subsystems are utilizing the network redundancy mechanism, while the other 2/3 are not. We call these two groups of storage subsystems *dual paths* systems and *single path* systems, respectively. In our data set, there are very few disk models used in both configurations; other disk models are mainly used in either *dual paths* systems or *single path* systems. Therefore, we cannot further break down the results based on disk models and shelf enclosure models.

Finding (7): Storage subsystems configured with network redundancy mechanisms experience much lower (30-40% lower) AFR than other systems. AFR for *physical interconnects* is reduced by 50-60%.

Implications: Network redundancy mechanisms such as multipathing can greatly improve the reliability of storage subsystems.

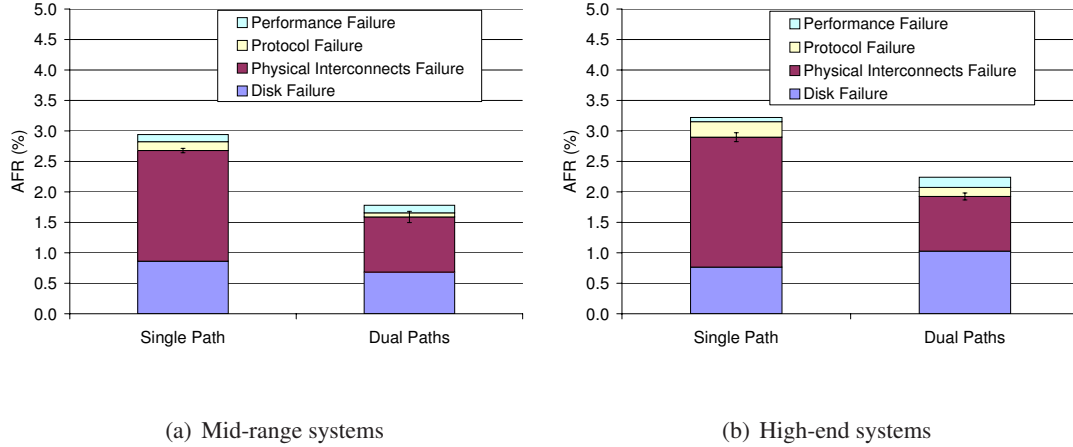


Figure 7. AFR for storage subsystems broken down by the number of paths. The error bars show 99.9% confidence intervals for physical interconnect failures.

Figure 7(a) and 7(b) show the AFR for storage subsystems in mid-range and high-end systems, respectively. As expected, secondary path reduces *physical interconnect failures* by 50-60% ($1.82 \pm 0.07\%$ versus $0.91 \pm 0.09\%$ and $2.13 \pm 0.07\%$ versus $0.90 \pm 0.06\%$), with little impact on other failure types. Since *physical interconnect failure* is just a subset of all *storage subsystem failures*, AFR for *storage subsystems* is reduced by 30-40%. This indicates that multipathing is an exceptionally good redundancy mechanism that delivers reduction of failure rates as promised. As we applied a T-test on these results, we found out that for both mid-range and high-end systems the observation is significant at the 99.9% confidence interval, indicating that the data strongly support the hypothesis that physical interconnect failures are reduced by multipathing configuration.

However, the observation also tells us that there is still further potential in network redundancy mechanism designs. For example, given that the probability for one network to fail is about 2%, the idealized probability for two networks to both fail should be a few magnitudes lower (about 0.04%). But the AFR we observe is far from the ideal number.

One reason is that not only failures from networks between shelves contribute to *physical interconnect failures*; other failures, such as shelf backplane errors, can also lead to *physical interconnect failures*, while multipathing does not provide redundancy for shelf backplane. Another possible reason is that most modern host adapters support more than one port, and each port can be used as a “logical” host adapter. If two independent networks are initiated by two “logical” host adapters sharing the same physical host adapter, a host adapter failure can cause failures of both networks.

4.4 Disk Positioning in the Shelf

As stated before, all the shelves we study in this report have 14 bays, and thus can accommodate at most 14 disks. However, a shelf doesn’t necessarily host 14 disks all the time, from the logs we analyzed, a shelf has on average 11.6 disks in it. Although not shown in Figure 1, there is an alarm in each shelf, located at one end of the shelf. It’s there as a means to notify the system administrator once some messages for users are generated.

Since disks are generally sensitive to vibration, there’s a concern that the vibration caused by alarms could potentially impact AFRs of all the disks in a shelf. A previous study shows that disk failure rate decreases

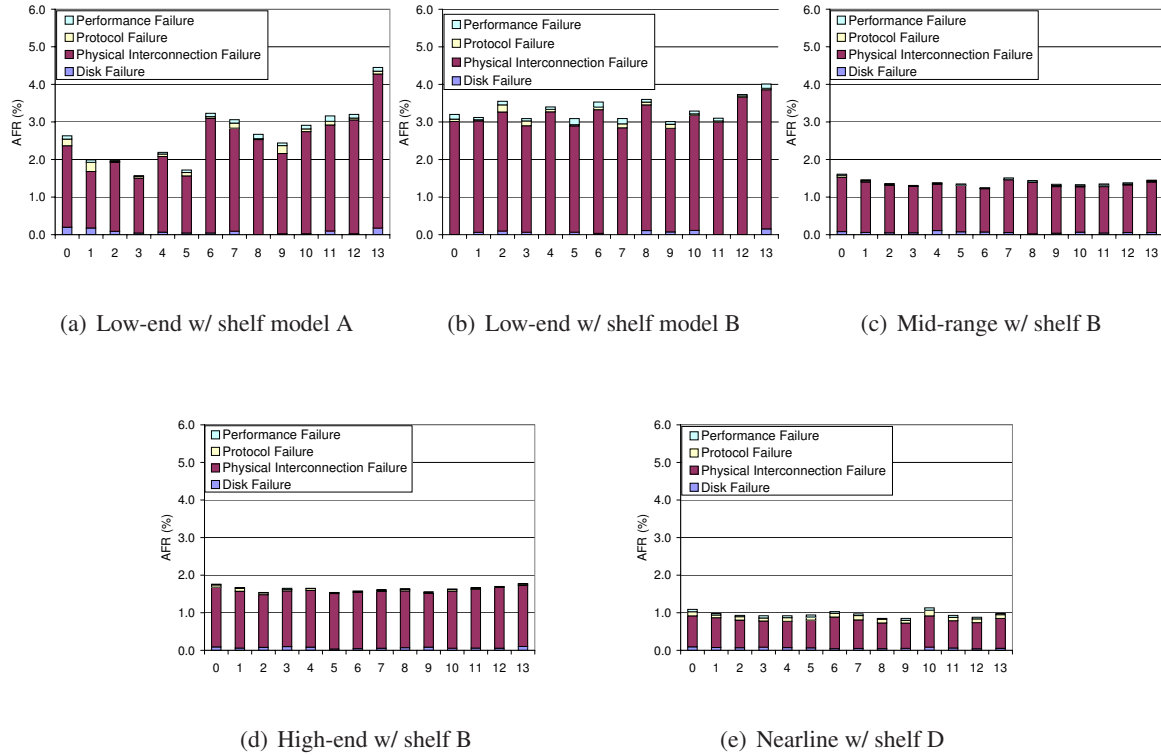


Figure 8. AFR for storage subsystems by bay id.

when bay id increases [12]. Combining the assumption that the alarm is at the location where bay id starts, that study support the hypothesis that alarm vibration might increase disk failure rate.

Figure 8 shows AFR for nearline secondary and low-end, mid-range and high-end primary storage systems, grouped by bay id. To avoid the results being skewed by the problematic disk family, we exclude data from storage subsystems using *Disk H*. From Figure 8(b) to 8(e), disk AFR and bay id doesn't show obvious correlation. However, as Figure 8(a) shows, the combination of shelf enclosure model and bay id show correlation with disk AFR. For shelf enclosure model A, disk AFR increases when bay id increases. Given that shelf model A is an internal shelf, which means that it is embedded in the storage system node, a possible reason for the correlation is that the end of the shelf with the largest bay id is closer to the heat flow in the node, and excessive heat increases AFR of disks in bays close to that end.

5 Statistical Properties of Storage Subsystem Failures

An important aspect of storage subsystem failures is their statistical properties. Understanding the statistical properties such as failure distribution of modern storage subsystems is necessary to build right testbed and fault injection models to evaluate existing resiliency mechanisms and to develop better ones. For example, some researchers have assumed a constant failure rate, which means an exponentially distributed time between failures, and that failures are independent, when calculating the expected time to failure for a RAID [14].

Figure 9 illustrates how disks are laid out in storage subsystems. As Figure 9 shows, multiple disks are mounted in one shelf enclosure and share the cooling service, power supply, and intrashelf connectivity provided by the shelf enclosure.

