

# Improving Reliability of Energy-Efficient Parallel Storage Systems by Disk Swapping

Shu Yin, Xiaojun Ruan, Adam Manzanares, Zhiyang Ding, Jiong Xie, James Majors, and Xiao Qin  
Department of Computer Science and Software Engineering

Auburn University, Auburn, AL 36849

Email: {szy0004, xzr0001, acm0008, dingzhi, jzx0009, majorjh, xqin}@auburn.edu

**Abstract**—The Popular Disk Concentration (PDC) technique and the Massive Array of Idle Disks (MAID) technique are two effective energy saving schemes for parallel disk systems. The goal of PDC and MAID is to skew I/O load towards a few disks so that other disks can be transitioned to low power states to conserve energy. I/O load skewing techniques like PDC and MAID inherently affect reliability of parallel disks because disks storing popular data tend to have high failure rates than disks storing cold data. To achieve good tradeoffs between energy efficiency and disk reliability, we first present a reliability model to quantitatively study the reliability of energy-efficient parallel disk systems equipped with the PDC and MAID schemes. Then, we propose a novel strategy—disk swapping—to improve disk reliability by alternating disks storing hot data with disks holding cold data. We demonstrate that our disk-swapping strategies not only can increase the lifetime of cache disks in MAID-based parallel disk systems, but also can improve reliability of PDC-based parallel disk systems.

**Keywords**—Parallel disk system, energy conservation, reliability, load balancing

## I. INTRODUCTION

Parallel disk systems, providing high-performance data-processing capacity, are of great value to large-scale parallel computers [1]. A parallel disk system comprised of an array of independent disks can be built from low-cost commodity hardware components. In the past few decades, parallel disk systems have increasingly become popular for data-intensive applications running on massively parallel computing platforms [2].

Existing energy conservation techniques can yield significant energy savings in disks. While several energy conservation schemes like cache-based energy saving approaches normally have marginal impact on disk reliability, many energy-saving schemes (e.g., dynamic power management and workload skew techniques) inevitably have noticeable adverse impacts on storage systems [3][4]. For example, dynamic power management (DPM) techniques save energy by using frequent disk spin-downs and spin-ups, which in turn can shorten disk lifetime [5] [6] [7], redundancy techniques [8] [9] [10] [11], workload skew [12] [13] [14], and multi-speed settings [15] [16]. Unlike DPM, workload-skew techniques such as MAID [17] and PDC [18] move popular data sets to a subset of disks arrays acting as workhorses, which are kept busy in a way that other disks can be turned into the standby mode to save energy. Compared with disks storing cold data, disks archiving hot data inherently have higher risk of breaking down.

Unfortunately, it is often difficult for storage researchers to improve reliability of energy-efficient disk systems. One of the main reasons lies in the challenge that every disk energy-saving research faces today, how to evaluate reliability impacts of power management strategies on disk systems. Although reliability of disk systems can be estimated by simulating the behaviors of energy-saving algorithms, there is lack of fast and accurate methodology to evaluate reliability of modern storage systems with high-energy efficiency. To address

this problem, we developed a mathematical reliability model called MINT to estimate the reliability of a parallel disk system that employs a variety of reliability-affecting energy conservation techniques [19].

In this paper, we first study the reliability of a parallel disk system equipped with two well-known energy-saving schemes—the PDC [18] and the MAID [17] technique. Preliminary results show that the reliability of PDC is slightly higher than that of MAID under light workload. We also observed that MAID is noticeably more reliable than PDC with relatively high data-access rates.

Since PDC does not support any data replication, the failure of a disk can cause the failure of the entire parallel disk system. Furthermore, I/O load skewing techniques like PDC and MAID inherently affect reliability of parallel disks because of two reasons: First, disks storing popular data tend to have high I/O utilization than disks storing cold data. Second, disks with higher utilization are likely to have higher risk of breaking down. To address the adverse impact of load skewing techniques on disk reliability, a disk swapping strategy was proposed to improve disk reliability in PDC by alternating disks storing hot data with disks holding cold data. Additionally, the disk swapping scheme was applied to MAID by switching the roles of data disks and cache disks. We evaluate impacts of the disk swapping scheme on the reliability of MAID-based parallel disk systems.

In this paper, our contributions are as follows:

- 1) We studied two reliability model for Popular Data Concentration scheme (PDC) and Massive Array of Idle Disks (MAID) based on Mathematical Reliability Models for Energy-efficient Parallel Disk System (MINT) [19];
- 2) We built a disk swapping mechanism to improve reliability of various load skewing techniques.
- 3) We studied the impacts of the disk swapping schemes on the reliability of PDC and MAID.

The remainder of this paper is organized as follows. Section II presents the framework of the MINT model. In Section III, the disk swapping mechanism equipped with MAID is presented. Section IV studies the reliability impacts of disk swapping on PDC. Section V presents experimental results and performance evaluation. In Section VI, the related work is discussed. Finally, Section VII concludes the paper with discussions.

## II. MINT: A RELIABILITY MODELING FRAMEWORK

### A. Overview

MINT is a framework developed to model reliability of parallel disk systems employing energy conservation techniques [19]. In the MINT framework, we studied the reliability impacts of two well-known energy-saving techniques - the Popular Disk Concentration technique (PDC) and the Massive Array of Idle Disks (MAID). One critical module in MINT is to model how PDC and MAID affect

the utilization and power-state transition frequency of each disk in a parallel disk system. Another important module developed in MINT is to calculate the annual failure rate of each disk as a function of the disk’s utilization, power-state transition frequency as well as operating temperature. Given the annual failure rate of each disk in the parallel disk system, MINT is able to derive the reliability of an energy-efficient parallel disk system. As such, we used MINT to study the reliability of a parallel disk system equipped with the PDC and MAID techniques.

Fig. 1 outlines the MINT reliability modeling framework. MINT is composed of a single disk reliability model, a system-level reliability model, and three reliability-affecting factors—temperature, power state transition frequency (hereinafter referred to as transition frequency or frequency) and utilization. Many energy-saving schemes (e.g., PDC [18] and MAID [17]) inherently affect reliability-related factors like disk utilization and transition frequency. Given an energy optimization mechanism, MINT first transfers data access patterns into the two reliability-affecting factors—frequency and utilization. The single disk reliability model can derive individual disk’s annual failure rate from utilization, power-state transition frequency, age, and temperature because these parameters are key reliability-affecting factors. Each disk’s reliability is used as input to the system-level reliability model that estimates the annual failure rate of parallel disk systems. For simplicity without losing generality, we considered in MINT four reliability-related factors, namely: disk utilization, age, temperature, and power-state transitions. This assumption does not necessarily indicate by any means that there are only four parameters affecting disk reliability. Other factors having impacts on reliability include: handling, humidity, voltage variation, vintage, duty cycle, and altitude [20]. That means if a new factor has to be taken into account, one can extend the single reliability model (see Section II-E) by integrating the new factor with other reliability-affecting factors in MINT. Since the infant mortality phenomenon is out the scope of this study, we pay attention to disks that are more than one year old.

### B. Disk Utilization

Disk utilization, a reliability-related factor, can be characterized as the fraction of active time of a disk drive out of its total powered-on-time [21]. In our single disk reliability model, the impacts of disk utilization on reliability is good way of providing a baseline characterization of disk annual failure rate (AFR). Pinheiro *et al.* studied the impact of utilization on AFR across different disk age groups [21]. They categorized disk utilization in three levels—low, medium, and high. Since the single disk reliability model needs a baseline AFR derived from a numerical value of utilization, we applied the polynomial curve-fitting technique to model the baseline value of a single disk’s AFR as a function of utilization. Thus, the baseline value (i.e., *BaseValue* in Eq. 1) of AFR for a disk can be calculated from the disk’s utilization.

### C. Temperature

Temperature is often considered as the most important environmental factor affecting disk reliability. For example, results from Google show that at very high temperatures, higher failure rates are associated with higher temperatures. In the low and middle temperature ranges, failure rate decreases when temperature increases [21].

In the MINT model, the temperature factor is a multiplier to base failure rates, which reflect reliability at base environmental conditions (see, for example, [20]). The temperature factor (i.e., *TemperatureFactor* in Eq. 1) is set to 1 when temperature is 25°C because room temperatures of many data centers are kept to 25°C

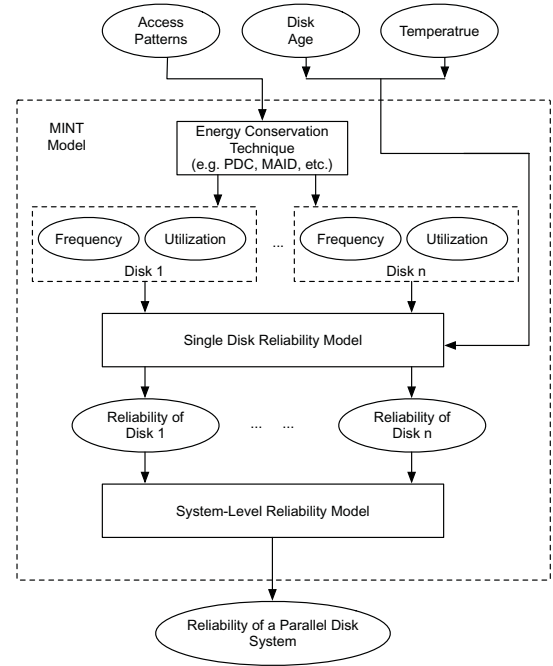


Fig. 1: Overview of the MINT reliability modeling methodology

by cooling systems. Suppose  $T$  is the average temperature, we define the temperature factor in case of  $T$  as  $T/25$  if  $T$  is larger than 25°C. When  $T$  exceeds 45°C, the temperature factor becomes a constant (i.e.,  $1.8 = 45/25$ ) because the cooling systems won’t let the room temperature higher than that.

### D. Power-State Transition Frequency

To conserve energy, power management policies turn idle disks from the active state into standby. The disk power-state transition frequency (or frequency for short) is often measured as the number of power-state transitions (i.e., from active to standby or vice versa) per month. The reliability of an individual disk is affected by power-state transitions and, therefore, the increase in failure rate as a function of power-state transition frequency has to be added to a baseline failure rate (see Eq. 1 in the next subsection).

### E. Single-Disk Reliability Model

Single-disk reliability can not be accurately described by one valued parameter because the disk drive reliability is affected by multiple factors (see Sections II-B, II-C, and II-D). We first compute a baseline failure rate as a function of disk utilization. Secondly, the temperature factor is used as a multiplier to the baseline failure rate. Finally, we add frequency to the baseline value of the annual failure rate. Hence, the failure rate  $R$  of an individual disk can be expressed as:

$$R = \alpha \times BaseValue \times TemperatureFactor + \beta \times FrequencyAdder \quad (1)$$

where *BaseValue* is the baseline failure rate derived from disk utilization (see Section II-B), *TemperatureFactor* is the temperature multiplier (see Section II-C), *FrequencyAdder* is the power-state transition frequency adder to the baseline failure rate (see Section II-D), and  $\alpha$  and  $\beta$  are two coefficients to reliability  $R$ . If reliability  $R$  is more sensitive to frequency than to utilization

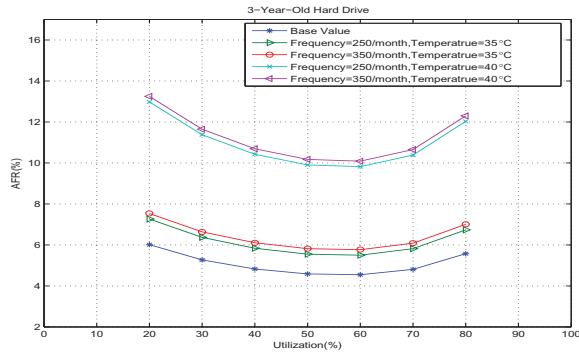


Fig. 2: Impacts of Combined Factors on the Annual Failure Rate of a 3-Year-Old HDD (Single Disk Reliability Model)

and temperature, then  $\beta$  must be greater than  $\alpha$ . Otherwise,  $\beta$  is smaller than  $\alpha$ . In either cases,  $\alpha$  and  $\beta$  can be set in accordance with  $R$ 's sensitivities to utilization, temperature, and frequency. In our experiments, we assume that all the three reliability-related factors are equally important (i.e.,  $\alpha=\beta=1$ ). Ideally, extensive field tests allow us to analyze and test the two coefficients. Although  $\alpha$  and  $\beta$  are not fully evaluated by field testing, reliability results are valid because of two reasons: first, we have used the same values of  $\alpha$  and  $\beta$  to evaluate impacts of the two energy-saving schemes on disk reliability (see Section III); second, the failure-rate trend of a disk when  $\alpha$  and  $\beta$  are set to 1 are very similar to those of the same disk when the values of  $\alpha$  and  $\beta$  do not equal to 1.

With Eq. 1 in place, we can analyze a disk's reliability in turns of annual failure rate or AFR. Fig. 2 shows AFR of a three-year-old disk when its utilization is in the range between 20% and 80%. We observe from Fig. 2 that increasing temperature from 35°C to 40°C gives rise to a significant increase in AFR. Unlike temperature, power-state transition frequency in the range of a few hundreds per month has marginal impact on AFR. It is expected that when transition frequency is extremely high, AFR becomes more sensitive to frequency than to temperature.

### III. DISK SWAPPING IN MAID

#### A. MAID - Massive Arrays of Idle Disks

The MAID (Massive Arrays of Idle Disks) technique - developed by Colarelli and Grunwald - aims to reduce energy consumption of large disk arrays while maintaining acceptable I/O performance [17]. MAID relies on data temporal locality to place replicas of active files on a subset of cache disks, thereby allowing other disks to spin down. Fig. 3 shows that MAID maintains two types of disks - cache disks and data disks. Frequently accessed files are copied from data disks into cache disks, where the LRU policy is implemented to manage data replacement in cache disks. Replaced data is discarded by a cache disk if the data is clean; dirty data has to be written back to the corresponding data disk. To prevent cache disk from being overloaded, MAID can avoid copying data to cache disks that have reached their maximum bandwidth. Three parameters will be used in systems: (1) power management policy, by using which drives that have not seen any requests for a specified period are spun down to sleep, or an adaptive spin-down to active; (2) data layout, which is either linear, with successive blocks being placed on the same drive, or striped across multiple drives; (3) cache, which indicates the number of drives of the array which will be used for cache [17].

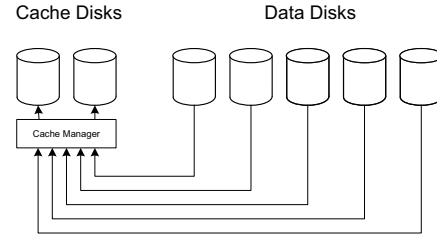


Fig. 3: The MAID System Structure

#### B. Improving Reliability of Cache Disks in MAID via Disk Swapping

Cache disks in MAID are more likely to fail than data disks due to the two reasons. First, cache disks are always kept active to maintain short I/O response times. Second, the utilization of cache disks is expected to be much higher than that of data disks. From the aspect of data loss, the reliability of MAID relies on the failure rate of data disks rather than that of cache disks. However, cache disks tend to be a single point of failure in MAID, which if the cache disks fail, will stop MAID from conserving energy. In addition, frequently replacing failed cache disks can increase hardware and management costs in MAID. To address this single point of failure issue and make MAID cost-effective, we designed a disk swapping strategy for enhancing the reliability of cache disks in MAID.

Fig. 4 shows the basic idea of the disk swapping mechanism, according to which disks rotate to perform the cache-disk functionality. In other words, the roles of cache disks and data disks will be periodically switched in a way that all the disks in MAID have equal chance to perform the role of caching popular data. For example, the two cache disks on the left-hand side in Fig. 4 are swapped with the two data disks on the right-hand side after a certain period of time (see Section V-C for circumstances under which disks should be swapped). For simplicity without losing generality, we assume that all the data disks in MAID initially are identical in terms of reliability. This assumption is reasonable because when a MAID system is built, all the new disks with the same model come from the same vendor. Initially, the two cache disks in Fig. 4 can be swapped with any data disk. After the initial phase of disk swapping, the cache disks are switched their role of storing replica data with the data disks with the lowest annual failure rate. In doing so, we ensure that cache disks are the most reliable ones among all the disks in MAID after each disk swapping process. It is worth noting that the goal of disk swapping is not to increase mean time to data loss, but is to boost mean time to cache-disk failure by balancing failure rates across all disks in MAID.

Disk swapping is advantageous to MAID, and the reason is two-fold. First, disk swapping further improves the energy efficiency of MAID because any failed cache disk can prevent MAID from effectively saving energy. Second, disk swapping reduces maintenance cost of MAID by making cache disks less likely to fail.

### IV. DISK SWAPPING IN PDC

#### A. PDC - Popular Data Concentration

The popular data concentration technique or PDC proposed by Pinheiro and Bianchini migrates frequently accessed data to a subset of disks in a disk array [18]. Fig. 5 demonstrates the basic idea behind PDC: the most popular files are stored in the far left disk, while the least popular files are stored in the far right disk. PDC relies on file popularity and migration to conserve energy in disk arrays, because

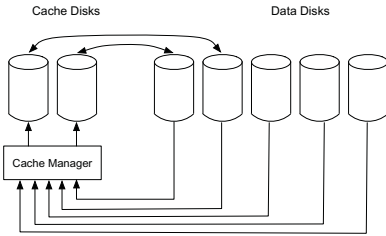


Fig. 4: Disk Swapping in MAID: The two cache disks on the left-hand side are swapped with the two data disks on the right-hand side.

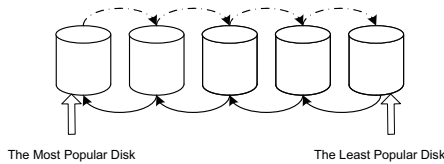


Fig. 5: PDC System Structure

several network servers exhibit I/O loads with highly skewed data access patterns. The migrations of popular files to a subset of disks can skew disk I/O load towards this subset, offering other disks more opportunities to be switched to standby to conserve energy. To avoid performance degradation of disks storing popular data, PDC aims at migrating data into a disk until its load is approaching the maximum bandwidth of the disk.

The main difference between MAID and PDC is that MAID makes data replicas on cache disks, whereas PDC lays data out across disk arrays without generating any replicas. If one of the cache disks fails in MAID, files residing in the failed cache disks can be found in the corresponding data disks. In contrast, any failed disk in PDC can inevitably lead to data loss. Although PDC tends to have lower reliability than MAID, PDC does not need to trade disk capacity for improved energy efficiency and I/O performance.

### B. Improving Reliability of PDC via Disk Swapping

Compared with disks storing popular files (hereinafter referred to as popular disks), disks archiving non-popular files tend to have lower failure rates because the utilization of popular disks is significantly higher than that of non-popular disks. In other words, popular disks have high chance of staying active serving I/O requests that access popular files, whereas non-popular disks are likely to be placed in the standby mode to conserve energy. Frequently repairing failed popular disks results in the following two adverse impacts. First, replacing failed hard drives with new ones can increase the overall hardware maintenance cost in PDC. Second, restoring disk data objects incurs extra energy overhead, thereby reducing the energy efficiency of PDC. To address these two problems of frequent repairing popular disks, we developed two disk-swapping strategies to improve the reliability of popular disks in PDC.

Fig. 6 shows the first disk-swapping strategy tailored for PDC. In this strategy called PDC-Swap1, the most popular disk is swapped with the least popular disk; the second most popular one is swapped with the second least popular disk; and the  $i$ th most popular disk is swapped with the  $i$ th least popular disk. Although the PDC-Swap1 is a straightforward approach, it does not necessarily lead to a balanced failure rates across all the disks in PDC. This is mainly because the graphs of AFR-utilization functions (see, for example,

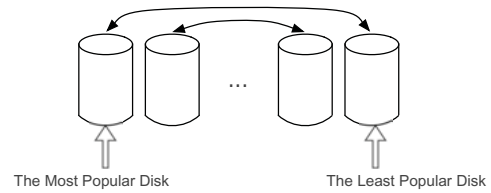


Fig. 6: PDC-Swap1- the first disk-swapping strategy in PDC: the  $i$ th most popular disk is swapped with the  $i$ th least popular disk.

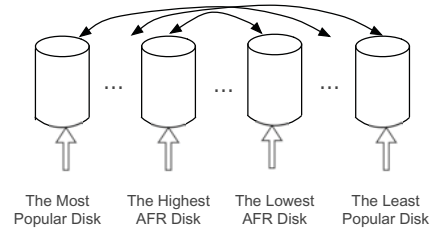


Fig. 7: PDC-Swap2 - The second disk-swapping strategy in PDC: the  $i$ th most popular disk is swapped with another disk with the  $i$ th lowest annual failure rate.

Fig. 2) are U-shaped curves, indicating that low disk utilization does not necessarily result in low annual failure rate (AFR). A potential disadvantage of the PDC-Swap1 strategy lies in the possibility that PDC-Swap1 is unable to noticeably improve the reliability of PDC. Such a disadvantage may occur when the utilization of the non-popular disks is so low that their AFRs are as high as those of popular disks.

To solve the above potential problems of PDC-Swap1, we proposed a second disk-swapping strategy in PDC. Fig. 7 outlines the basic idea of the second disk-swapping strategy called PDC-Swap2, which attempts to switch the most popular disk with the one with the lowest AFR. In general, PDC-Swap2 switches the  $i$ th most popular disk with another disk with the  $i$ th lowest AFR. The PDC-Swap2 strategy guarantees that after each swapping procedure, disks with low failure rates start serving as popular disks. Unlike PDC-Swap1, PDC-Swap2 can substantially improve the reliability of PDC by balancing the failure rates across all the popular and non-popular disks.

## V. EXPERIMENTAL RESULTS

### A. Experimental Setup

We developed a simulator to validate the reliability models for MAID and PDC. It might be unfair to compare the reliability of MAID and PDC using the same number of disks, since MAID trades extra cache disks for high energy efficiency. To make fair comparison, we considered two system configurations for MAID. The first configuration referred to as MAID-1 employs existing disks in a parallel disk system as cache disks to store frequently accessed data. Thus, the first configuration of MAID improves energy efficiency of the parallel disk system at the cost of capacity. In contrast, the second configuration— called MAID-2—needs extra disks to be added to the disk system to serve as cache disks.

Our experiments were started by evaluating the reliability of PDC as well as MAID-1 and MAID-2. Then, we studied the reliability impacts of the proposed disk-swapping strategies on both PDC and MAID. We simulated PDC, MAID-1, and MAID-2 along with the



disk-swapping strategies in two parallel disk systems described in Table I. For the MAID-1 configuration, there are 5 cache disks and 15 data disks. In the disk system for the MAID-2 configuration, there are 5 cache disks and 20 data disks. As for the case of PDC, we fixed the number of disks to 20. Thus, we studied MAID-2 and PDC using a parallel disk system with 20 disks; we used a similar disk system with totally 25 disks to investigate MAID-1. We varied the file access rate in the range between 0 to  $10^6$  times per month. The average file size considered in our experiments is 300KB. The base operating temperature is set to  $35^\circ\text{C}$ . In this study, we focused on read-only workload. Nevertheless, the MINT model should be readily extended to capture the characteristics of read/write workloads.

TABLE I: The characteristics of the simulated parallel disk system used to evaluate the reliability of PDC, MAID-1, and MAID-2.

Energy-efficiency Scheme	Number of Disks	File Access Rate (No. per month)	File Size (KB)
PDC	20 data (20 in total)	$0\sim 10^6$	300
MAID-1	15 data+5 cache (20 in total)	$0\sim 10^6$	300
MAID-2	20 data+5 cache (25 in total)	$0\sim 10^6$	300

### B. PDC and MAID without Disk Swapping

Let us first examine the reliability of PDC, MAID-1, and MAID-2 in which the disk-swapping strategies are not incorporated. Fig. 8 shows the utilization as a function of file access rate. Fig. 8 indicates that PDC's utilization is faster than those of MAID-1 and MAID-2, since the utilization of PDC quickly approaches to 90%. The main reason is that under dynamic I/O workload conditions, PDC needs to spend time in migrating data between popular and non-popular disks. Unlike PDC, MAID is not very sensitive to dynamically changing workloads. Interestingly, the utilization of MAID-1 grows faster than that of MAID-2, because MAID-2 has five more disks compared to MAID-1.

Recall that disk utilization is one important reliability-affecting factor. Fig. 9 shows that the annual failure rates (AFR) of PDC, MAID-1, and MAID-2 are changing in accordance with disk utilization. We observe from Fig. 9 that the AFR value of PDC keeps increasing from 5.6% to 8.3% when the file access rate is larger than  $15 \times 10^3$  times/month. We attribute this trend to high disk utilization due to data migrations. More interestingly, if the file access rate is lower than  $15 \times 10^3$ , AFR of PDC slightly reduces from 5.9% to 5.6% when the access rate is increased from  $5 \times 10^3$  to  $1.5 \times 10^4$ . This result can be explained by the nature of the utilization function that is concave rather than linear. The concave nature of the utilization function is consistent with the empirical results reported in [21]. When the file access rate  $1.5 \times 10^4$ , the disk utilization is approximately 50%, which is the turning point of the AFR-utilization function.

The AFR value of MAID continues decreasing from 6.3% to 5.8% with the increasing file access rate. This declining trend can be explained by two reasons. First, increasing the file access rates reduces the number of power-state transitions. Second, the disk utilization is close to 40%, which is in the declining part of the AFR-utilization function. When the access rate keeps increasing to  $10^6$  times/month, one important observation from Fig. 9 is that when access rate is higher than  $7 \times 10^5$ , AFR of MAID-1 is getting higher than that of MAID-2. This trend is reasonable because MAID-1's disk utilization keeps rising up over 60% (see Fig. 8) when access rate is higher than  $7 \times 10^5$ . According to Fig. 2, AFRs stop rising

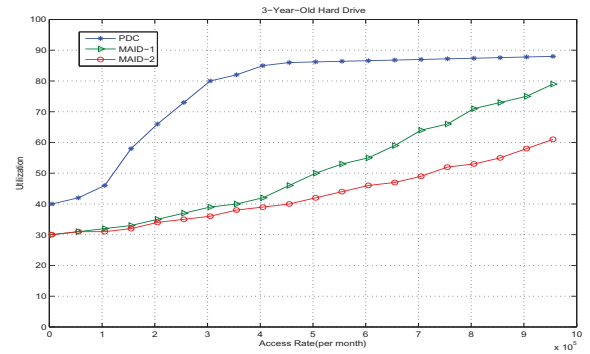


Fig. 8: Utilization Comparison of the PDC and MAID Access Rate Impacts on Utilization

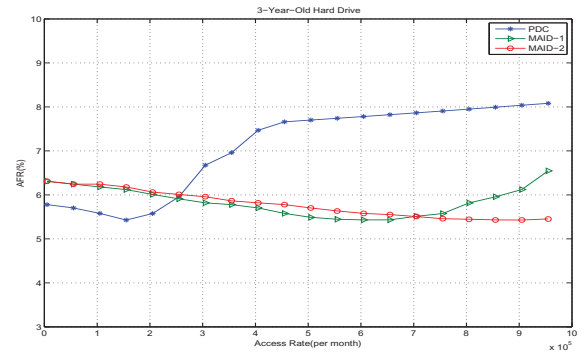


Fig. 9: AFR Comparison of the PDC and MAID Access Rate Impacts on AFR(Temperature= $35^\circ\text{C}$ )

after utilization becomes higher than 60%. Hence, we can conclude that after access rate hit  $9 \times 10^5$  times/month, AFR of MAID-2 is expected to stop increasing.

### C. Preliminary Results of the Disk-Swapping Strategies

A key issue of the disk-swapping strategies is to determine circumstances under which disks should be swapped in order to improve disk system reliability. One straightforward way to address this issue is to periodically initiate the disk-swapping process. For example, we can swap disks in MAID and PDC once every year. Periodically swapping disks, however, might not always enhance the reliability of parallel disk systems. For instance, swapping disks under very light workloads cannot substantially improve disk system reliability. In some extreme cases, swapping disks under light workload may worsen disk reliability due to overhead of swapping. As such, our disk-swapping strategies do not periodically swap disks. Rather, the disk-swapping process is initiated when the average I/O access rates exceed a threshold. In our experiments, we evaluated the impact of this access-rate threshold on the reliability of a parallel disk system. More specifically, the threshold is set to  $2 \times 10^5$ ,  $5 \times 10^5$ , and  $8 \times 10^5$  times/month, respectively. These three values are representative values for the threshold because when the access rate hits  $5 \times 10^5$ , the disk utilization lies in the range between 80% and 90% (see Fig. 8), which in turn ensures that AFR increases with the increasing value of utilization (see Fig. 2).

Figs. 10, 11, and 12 reveal the annual failure rates (AFR) of MAID-1 and MAID-2 with and without using the proposed disk-swapping

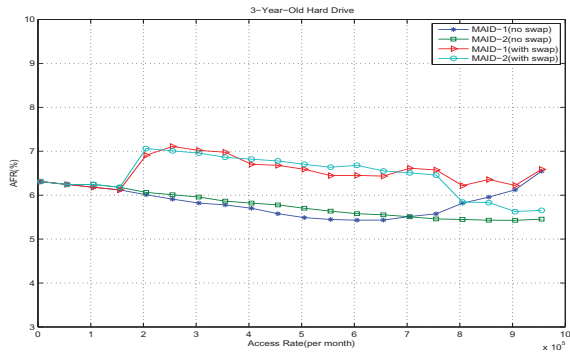


Fig. 10: Utilization Comparison of the MAID  
Access Rate Impacts on AFR (Threshold =  $2 * 10^5$  No./month)

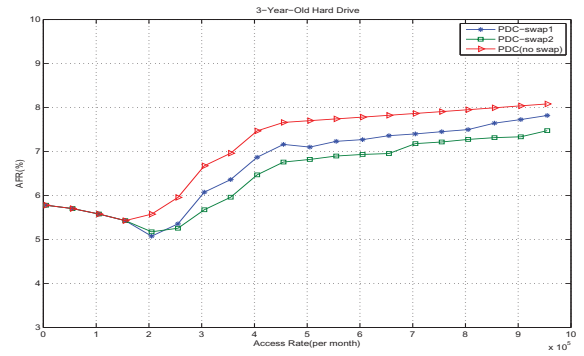


Fig. 13: Utilization Comparison of the PDC  
Access Rate Impacts on AFR (Threshold =  $2 * 10^5$  No./month)

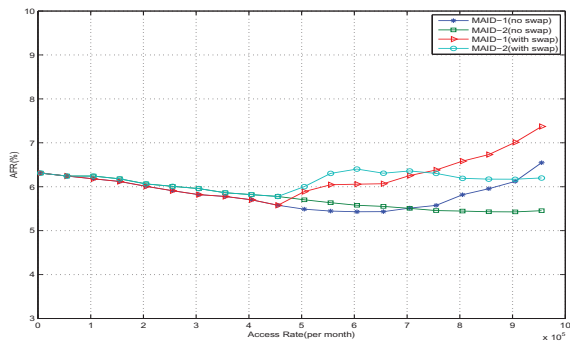


Fig. 11: Utilization Comparison of the MAID  
Access Rate Impacts on AFR (Threshold =  $5 * 10^5$  No./month)

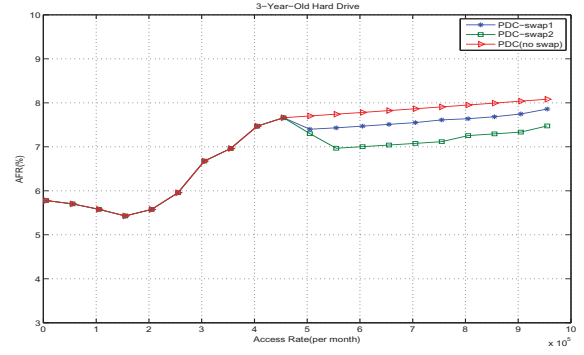


Fig. 14: Utilization Comparison of the PDC  
Access Rate Impacts on AFR (Threshold =  $5 * 10^5$  No./month)

strategy. The results plotted in Figs. 10, 11, and 12 show that for both MAID-1 and MAID-2, the disk-swapping process reduces the reliability of data disks in the disk system. We attribute the reliability degradation to the following reasons. MAID-1 and MAID-2 only store replicas of popular data; the reliability of the entire disk system is not affected by failures of cache disks. The disk-swapping processes increase the average utilization of data disks, thereby increasing the AFR values of data disks. Nevertheless, the disk-swapping strategy has its own unique advantage. Disk swapping is intended to reduce hardware maintenance cost by increasing the

lifetime of cache disks. In other words, disk swapping is capable of extending the Mean Time To Failure or MTTF [21] of the cache disks.

We also observed from Figs. 10, 11, and 12 that for the MAID-based disk system with the disk-swapping strategy, a small threshold leads to a low AFR. Compared with the other two thresholds, the  $2 * 10^5$  threshold results in the lower AFR. The reason is that when the access rate is  $2 * 10^5$  No./month, the disk utilization is around 35% (see Fig.8), which lies in the monotone decreasing area of the curve shown in Fig. 2. Thus, disk swapping reduces AFR for a while until the disk utilization reaches 60%.

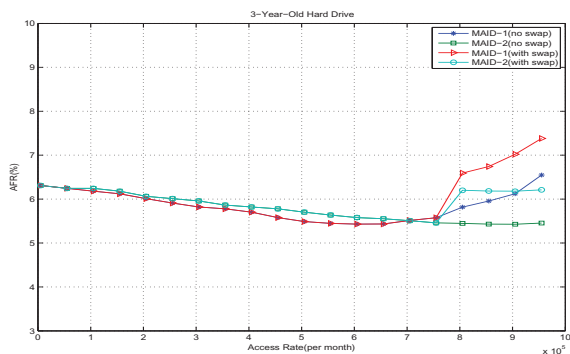


Fig. 12: Utilization Comparison of the MAID  
Access Rate Impacts on AFR (Threshold =  $8 * 10^5$  No./month)

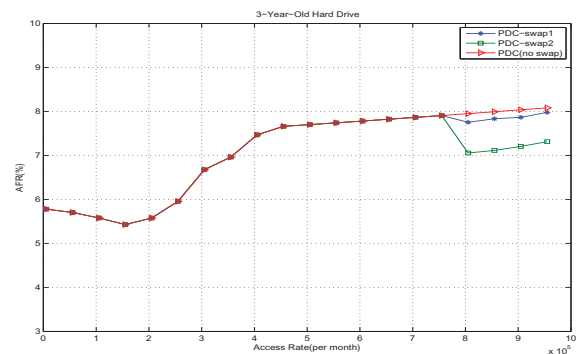


Fig. 15: Utilization Comparison of the PDC  
Access Rate Impacts on AFR (Threshold =  $8 * 10^5$  No./month)

Figs. 13, 14, and 15 show the AFR values of the parallel disk system where the PDC-Swap1 and PDC-Swap2 strategies are employed. For comparison purpose, we also examined AFRs of the same disk without using any disk-swapping strategy. The results plotted in Figs. 13, 14, and 15 indicate that the PDC-Swap1 and PDC-Swap2 strategies noticeably improve the reliability of the PDC-based parallel disk system by reducing the annual failure rate of the system. After comparing the three PDC-based disk systems, we observed that disk swapping reduces AFRs by nearly 5%. Compared with PDC-Swap1 strategy, PDC-Swap2 improves the reliability of the PDC-based system in a more efficient way. PDC-Swap2 is better than PDC-Swap1 because of the following reasons. First, PDC-Swap1 simply switches the most popular disks with the least popular disks whereas PDC-Swap2 swaps the most popular disks with the ones that have the least AFR values (see Section V-C). Second, low utilization does not necessarily lead to a low AFR value due to the U-shaped curve appeared in Fig. 2. Third, the least popular disks may not be the ones with the lowest AFR value, meaning that PDC-Swap1 cannot guarantee that the popular disks are always swapped with the disks with the lowest AFRs.

After comparing the results presented in Figs. 13, 14, and 15, we observe that the improvement in reliability is sensitive to the threshold. More specifically, an increased threshold can give rise to the increase in the reliability improvement. This phenomenon can be explained as follows. After the disk swapping process, the utilization of the most popular disk is reduced down to one with the lowest AFR of the entire parallel disk system. A higher threshold implies a larger utilization discrepancy between the pair of swapped disks. Importantly, regardless of the threshold value, AFR of the disk system continues increasing at a slow pace with the increasing value of access rate. This trend indicates that after each disk swapping, the utilizations of those disks with low AFRs are likely to be kept at a high level, which in turn leads to an increasing AFR of the entire disk system.

## VI. RELATED WORK

A hard disk drive (HDD) is a complex dynamic system made up of various electrical, electronic, and mechanical components [22]. An array of techniques were developed to save energy in single HDDs. Energy dissipation in disk drives can be reduced at the I/O level (e.g., dynamic power management [23][7] and multi-speed disks [6]), the operating system level (e.g., power-aware caching/prefetching [9][16]), and the application level (e.g. software DMP [24] and cooperative I/O [25]). Existing energy-saving techniques for parallel disk systems often rely on one of the two basic ideas - power management and workload skew. Power management schemes conserve energy by turning disks into standby after a period of idle time. Although multi-speed disks are not widely adopted in storage systems, power management has been successfully extended to address the energy-saving issues in multi-speed disks [6][15][26]. The basic idea of workload skew is to concentrate I/O workloads from a large number of parallel disks into a small subset of disks allowing other disks to be placed in the standby mode [18][17][27][28].

Recent studies show that both power management and workload skew schemes inherently impose adverse impacts on disk systems [3][4]. For example, the power management schemes are likely to result in a huge number of disk spin-downs and spin-ups that can significantly reduce hard disk lifetime. The workload skew techniques dynamically migrates frequently accessed data to a subset of disks [29] [30], which inherently have higher risk of breaking

down than other disks usually being kept on standby. Disks that store popular data tend to have high failure rates due to extremely unbalanced workload. Thus, the popular data disks have a strong likelihood to become reliability bottleneck. The design of our MINT is orthogonal to the aforementioned energy saving studies, because MINT is focused on reliability impacts of the power management and workload skew schemes in parallel disks.

A malfunction of any components in a hard disk drive could lead to a failure of the disk. Reliability—one of the key characteristics of disks—can be measured in terms of mean-time-between-failure (MTBF). Disk manufacturers usually investigate MTBFs of disks either by laboratory testing or mathematical modeling. Although disk drive manufacturers claim that MTBF of most disks is more than 1 million hours [31], users have experienced a much lower MTBF from their field data [20]. More importantly, it is challenging to measure MTBF because of a wide range of contributing factors including disk age, utilization, temperature, and power-state transition frequency [20].

A handful of reliability models have been successfully developed for storage systems. For example, Paris *et. al* investigated an approach to computing both average failure rate and mean time to failure in distributed storage systems [32]; Elerath and Pecht proposed a flexible model for estimating reliability of RAID storage [33]; and Xin *et. al* developed a model to study disk infant mortality [34]. Unlike these reliability models tailored for conventional parallel and distributed disk systems, our MINT model pays special attention to reliability of parallel disk systems coupled with energy-saving mechanisms.

Very recently, Xie and Sun developed an empirical reliability model called PRESS (Predictor of Reliability for Energy Saving Schemes) [4]. The PRESS model can be used to estimate reliability of an entire disk array [4]. To fully leverage PRESS to study the reliability of disk arrays, one has to properly simulate the disk arrays. Our MINT approach differs itself from PRESS in the sense that the goal of MINT is to evaluate reliability of disk systems by modeling the behavior of parallel disks where energy conservation mechanisms are integrated.

Swapping mechanisms have been thoroughly studied in the arena of memory and file systems. For example, Paul *et. al* developed an efficient virtual memory swapping system - called LocalSwap - to improve performance of clusters [35]; Plank addressed the issue of checkpoint placement and its impact on the performance of the PVM platform [36]; Pei and Edward investigated the performance of a file system based on the LRU-SP(Least-Recently-Used with Swapping) policy [37]. Our disk swapping approaches are fundamentally different from the aforementioned swapping mechanisms in the sense that the goal of disk swapping is to improve the reliability of energy-efficient parallel disk systems by balancing the failure rates of parallel disks.

## VII. CONCLUSIONS

This paper presents a reliability model to quantitatively study the reliability of energy-efficient parallel disk systems equipped with the Massive Array of Idle Disks (MAID) technique and the Popular Disk Concentration (PDC) technique. Note that MAID and PDC are two effective energy-saving schemes for parallel disk systems. MAID and PDC aim to skew I/O load towards a few disks so that other disks can be transitioned to low power states to conserve energy. I/O load skewing techniques like MAID and PDC inherently affect reliability of parallel disks because disks storing popular data tend to have high failure rates than disks storing cold data. To address the reliability issue in MAID and PDC, we developed disk-swapping

strategies to improve disk reliability by alternating disks storing hot data with disks holding cold data. Additionally, we quantitatively evaluate the impacts of the disk-swapping strategies on reliability of MAID-based and PDC-based parallel disk systems. We demonstrate that the disk-swapping strategies not only can increase the lifetime of cache disks in MAID-based parallel disk systems, but also can significantly improve reliability of PDC-based parallel disk systems.

Future directions of this research can be performed in the following. First, we will extend the MINT model to investigate mixed read/write workloads in the future. Second, we will investigate a fundamental trade-off between reliability and energy-efficiency in the context of energy-efficient disk arrays. A tradeoff curve will be used as a unified framework to justify whether or not it is worth trading reliability for high energy efficiency. Third, we will study the most appropriate conditions under which disk-swapping processes should be initiated. Last, we will develop a multi-swapping mechanism that aims at balancing the utilization of each disk in a parallel disk system to maintain the failure rate the disk system at a low level.

#### ACKNOWLEDGMENT

The work reported in this paper was supported by the US National Science Foundation under Grants CCF-0845257 (CAREER), CNS-0757778 (CSR), CCF-0742187 (CPA), CNS-0917137 (CSR), CNS-0831502 (CyberTrust), CNS-0855251 (CRI), OCI-0753305 (CI-TEAM), DUE-0837341 (CCLI), and DUE-0830831 (SFS), as well as Auburn University under a startup grant and a gift (Number 2005-04-070) from the Intel Corporation.

#### REFERENCES

- [1] "The distributed-parallel storage system (dpss) home pages," <http://www.didc.lbl.gov/DPSS/>, June 2004.
- [2] P. Varman and R. Verma, "Tight bounds for prefetching and buffer management algorithms for parallel I/O systems," *IEEE Trans. Parallel Distrib. Syst.*, vol. 10, no. 12, pp. 1262–1275, 1999.
- [3] K. Bellam, A. Manzanares, X. Ruan, X. Qin, and Y.-M. Yang, "Improving reliability and energy efficiency of disk systems via utilization control," in *Proc. IEEE Symp. Computers and Comm.*, 2008.
- [4] T. Xie and Y. Sun, "Sacrificing reliability for energy saving: Is it worthwhile for disk arrays?" in *Proc. IEEE Symp. Parallel and Distr. Processing*, April 2008, pp. 1–12.
- [5] F. Douglass, P. Krishnan, and B. Marsh, "Thwarting the power-hungry disk," in *Proc. USENIX Winter 1994 Technical Conf.*, 1994, pp. 23–23.
- [6] D. Helmbold, D. Long, T. Sconyers, and B. Sherrord, "Adaptive disk spin—down for mobile computers," *Mob. Netw. Appl.*, vol. 5, no. 4, pp. 285–297, 2000.
- [7] K. Li, R. Kumpf, P. Horton, and T. Anderson, "A quantitative analysis of disk drive power management in portable computers," in *Proc. USENIX Winter Technical Conf.*, 1994, pp. 22–22.
- [8] E. Pinheiro, R. Bianchini, and C. Dubnicki, "Exploiting redundancy to conserve energy in storage systems," in *Proc. Joint Int'l Conf. Measurement and Modeling of Computer Systems*, 2006.
- [9] Q.-B. Zhu, F. David, C. Devaraj, Z.-M. Li, Y.-Y. Zhou, and P. Cao, "Reducing energy consumption of disk storage using power-aware cache management," in *Proc. Int'l Symp. High Performance Comp. Arch.*, Washington, DC, USA, 2004, p. 118.
- [10] J. Wang, H.-J. Zhu, and D. Li, "eraid: Conserving energy in conventional disk-based raid system," *IEEE Trans. Computers*, vol. 57, no. 3, pp. 359–374, 2008.
- [11] T. Xie, "Sea: A striping-based energy-aware strategy for data placement in raid-structured storage systems," *IEEE Trans. Computers*, vol. 57, no. 6, pp. 748–761, June 2008.
- [12] A. E. Papathanasiou and M. L. Scott, "Power-efficient server-class performance from arrays of laptop disks," 2004. [Online]. Available: <http://hdl.handle.net/1802/314>
- [13] S. Jin and A. Bestavros, "Gismo: A generator of internet streaming media objects and workloads." *ACM SIGMETRICS Performance Evaluation Review*, November 2001.
- [14] Q. Yang and Y.-M. Hu, "DCD - Disk Caching Disk: A new approach for boosting I/O performance," in *Proc. Int'l Symp. Computer Architecture*, May 1996, pp. 169–169.
- [15] S. Gurumurthi, A. Sivasubramaniam, M. Kandemir, and H. Franke, "Drpm: dynamic speed control for power management in server class disks," in *Proc. Int'l Symp. Computer Architecture*, June 2003, pp. 169–179.
- [16] S. Son and M. Kandemir, "Energy-aware data prefetching for multi-speed disks," in *Proc. Int'l Conf. Comp. Frontiers*, 2006.
- [17] D. Colarelli and D. Grunwald, "Massive arrays of idle disks for storage archives," in *Proc. ACM/IEEE Conf. Supercomputing*, 2002, pp. 1–11.
- [18] E. Pinheiro and R. Bianchini, "Energy conservation techniques for disk array-based servers," in *Proc. 18th Int'l Conf. Supercomputing*, 2004.
- [19] S. Yin, X. Ruan, A. Manzanares, and X. Qin, "How reliable are parallel disk systems when energy-saving schemes are involved?" in *Proc. IEEE International Conference on Cluster Computing (CLUSTER)*, 2009.
- [20] J. Elerath, "Specifying reliability in the disk drive industry: No more mtbf's," 2000, pp. 194–199.
- [21] E. Pinheiro, W.-D. Weber, and L. Barroso, "Failure trends in a large disk drive population," in *Proc. USENIX Conf. File and Storage Tech.*, February 2007.
- [22] J. Yang and F.-B. Sun, "A comprehensive review of hard-disk drive reliability," in *Proc. Annual Reliability and Maintainability Symp.*, 1999.
- [23] F. Douglass, P. Krishnan, and B. Marsh, "Thwarting the power-hungry disk," in *Proc. USENIX Winter 1994 Technical Conf.*, 1994, pp. 23–23.
- [24] S. W. Son, M. Kandemir, and A. Choudhary, "Software-directed disk power management for scientific applications," in *Proc. IEEE Int'l Parallel and Distr. Processing Symp.*, 2005.
- [25] A. Weissel, B. Beutel, and F. Belloso, "Cooperative I/O: a novel I/O semantics for energy-aware applications," in *Proc. the 5th Symp. Operating Systems Design and Implementation*. New York, NY, USA: ACM, 2002, pp. 117–129.
- [26] P. Krishnan, M. P. Long, and S. J. Vitter, "Adaptive disk spindown via optimal rent-to-buy in probabilistic environments," Durham, NC, USA, Tech. Rep., 1995.
- [27] X.-J. R. Run, A. Manzanares, S. Yin, Z.-L. Zong, and X. Qin, "Performance evaluation of energy-efficient parallel I/O systems with write buffer disks," in *Proc. 38th Int'l Conf. Parallel Processing*, Sept. 2009.
- [28] E. Pinheiro, R. Bianchini, E. Carrera, and T. Heath, "Load balancing and unbalancing for power and performance in cluster-based systems," *Proc. Workshop Compilers and Operating Sys. for Low Power*, September 2001.
- [29] X. J. Ruan, A. Manzanares, K. Bellam, Z. L. Zong, and X. Qin, "Daraw: A new write buffer to improve parallel I/O energy-efficiency," in *Proc. ACM Symp. Applied Computing*, 2009.
- [30] A. Manzanares, X. Ruan, S. Yin, and M. Nijim, "Energy-aware prefetching for parallel disk systems: Algorithms, models, and evaluation," *IEEE Int'l Symp. on Network Computing and Applications*, 2009.
- [31] B. Schroeder and G. Gibson, "Disk failures in the real world: what does an mtf of 1,000,000 hours mean to you?" in *Proc. USENIX Conf. File and Storage Tech.*, 2007, p. 1.
- [32] J.-F. Pâris, T. Schwarz, and D. Long, "Evaluating the reliability of storage systems," in *Proc. IEEE Int'l Symp. Reliable and Distr. Sys.*, 2006.
- [33] J. Elerath and M. Pecht, "Enhanced reliability modeling of raid storage systems," in *Proc. IEEE/IFIP Int'l Conf. Dependable Sys. and Networks*, 2007.
- [34] Q. Xin, J. Thomas, S. Schwarz, and E. Miller, "Disk infant mortality in large storage systems," in *Proc. IEEE Int'l Symp. Modeling, Analysis, and Simulation of Computer and Telecomm. Sys.*, 2005.
- [35] P. Werstein, X. Jia, and Z. Huang, "A remote memory swapping system for cluster computers," in *PDCAT '07: Proceedings of the Eighth International Conference on Parallel and Distributed Computing, Applications and Technologies*. Washington, DC, USA: IEEE Computer Society, 2007, pp. 75–81.
- [36] J. S. Plank, "Improving the performance of coordinated checkpoints on networks of workstations using raid techniques," in *SRDS '96: Proceedings of the 15th Symposium on Reliable Distributed Systems*. Washington, DC, USA: IEEE Computer Society, 1996, p. 76.
- [37] P. Cao, E. W. Felten, A. R. Karlin, and K. Li, "Implementation and performance of integrated application-controlled file caching, prefetching, and disk scheduling," *ACM Trans. Comput. Syst.*, vol. 14, no. 4, pp. 311–343, 1996.