Contents lists available at ScienceDirect

J. Parallel Distrib. Comput.

journal homepage: www.elsevier.com/locate/jpdc

Understanding the relationship between energy conservation and reliability in parallel disk arrays

Tao Xie^{a,*}, Yao Sun^b

^a San Diego State University, United States ^b Teradata Corporation, United States

ARTICLE INFO

Article history: Received 13 January 2010 Received in revised form 11 June 2010 Accepted 4 August 2010 Available online 13 August 2010

Keywords: Energy-aware systems Hardware reliability Parallel I/O Secondary storage

ABSTRACT

Energy conservation schemes based on power management or workload skew for disk arrays adversely affect disk reliability due to either workload concentration or frequent disk speed transitions. A thorough understanding of the relationship between energy-saving techniques and disk reliability is still an open problem, which prevents effective design of new energy-saving techniques and application of existing approaches in reliability-critical environments. This paper presents an empirical reliability model, called PRESS (Predictor of Reliability for Energy-Saving Schemes). Fed by operating temperature, disk utilization, and disk speed transition frequency, PRESS estimates the reliability of an entire disk array. Further, a new energy-saving strategy with reliability awareness named READ (Reliability and Energy Aware Distribution) is developed in the light of the insights provided by PRESS. Experimental results demonstrate that READ consistently performs better than existing approaches in performance and reliability while achieving a comparable level of energy consumption.

© 2010 Elsevier Inc. All rights reserved.

Journal of Parallel and Distributed Computing

1. Introduction

A hard disk drive (HDD) is a complex dynamic system that is made up of various electrical, electronic, and mechanical components [40]. A malfunction of any of these components could lead to a complete failure of a hard disk drive. While the capacity, spindle speed, form factor, and performance of hard disk drives have been enhanced rapidly, the reliability of hard disk drives is improving slowly. The primary reasons are that the hard disk manufacturing technology is constantly changing, and that the performance envelope of hard disk drives is incessantly pushed. MTBF (Mean Time Between Failure) and AFR (Annualized Failure Rate) are two related reliability metrics used by disk drive manufacturers, who have claimed that the MTBF of their enterprise products is more than 1 million hours, or roughly 114 years [30]. Storage system integrators and end users, however, have challenged the unrealistic reliability specification and have usually found a much lower MTBF from their field data [10]. Current reliability techniques are mainly leveraged on a variety of data redundancy mechanisms like data replication (RAID 1), paritybased protection (RAID 5), and Reed-Solomon erasure-correcting codes [25]. Still, maintaining a high level of reliability for a diskarray-based large-scale storage system with thousands of hard disk

drives is a major challenge because the very large number of disks dramatically lowers down the overall MTBF of the entire storage system [11]. The level of reliability for disk-array-based large-scale storage systems is far from satisfied [28].

Recently, energy conservation for disk arrays has been an important research topic in storage systems as they can consume 27% of the overall electricity in a data center [28]. A broad spectrum of technologies including power management [4,13,36,42], workload skew [5,24], RAID configuration [15,37], caching [3,43], and data placement [35,34] have been employed to save energy for disk arrays. However, power management based and workload skew based techniques (see Table 1), two typical energy saving schemes for disk arrays, negatively affect the lifetime of disks. Power management based energy conservation schemes like Multi-speed [4], DRPM [13], and Hibernator [42] frequently spin up or spin down disk drives, which obviously affects drives' lifetime. Workload skew oriented energy conservation techniques such as MAID [5] and PDC [26] utilize a subset of a disk array as a workhorse to store popular data so that other disks could have opportunities to have a rest to save energy. Apparently, very high disk utilization is detrimental for the reliability of those overly used disks, whose high failure rates degrade the reliability of the entire disk array.

Although most of the researchers who proposed the energysaving schemes above realized that their techniques could inherently affect disk reliability, only a few of them mentioned some intuitive ways such as limiting the power cycling of a disk to 10 times a day or rotating power-always-on disk role [37], to alleviate the side-effects of their schemes on disk reliability



^{*} Corresponding author. E-mail addresses: xie@cs.sdsu.edu (T. Xie), calvin.sun@teradata.com (Y. Sun).

^{0743-7315/\$ –} see front matter s 2010 Elsevier Inc. All rights reserved. doi:10.1016/j.jpdc.2010.08.006

 Table 1

 Energy conservation techniques.

Category	Example				
Power management	DRPM [13], Multi-speed [4], Hibernator [42], CMTPM, CMDRPM [36]				
Workload skew RAID configuration Caching Data placement	MAID [5], PDC [26] RAID Tuning [14], PARAID [37] PA-LRU, PB-LRU [43], CBSM [3] PDDL [34], PF+ [35]				

[37,43]. Still, a deep understanding of the relationship between energy-saving techniques and disk reliability is an open question.

To answer this question, a reliability model, which can quantify the impact of energy-saving-related reliability affecting factors (hereafter referred to as ESRRA factors) like operating temperature, speed transition frequency, and disk utilization on disk reliability, is fundamental. We present an empirical reliability model, called PRESS (Predictor of Reliability for Energy-Saving Schemes), which translates the three ESRRA factors into AFR (Annualized Failure Rate). The PRESS model provides us with a much needed understanding of the relationship between energy-saving and disk reliability. In the light of the insights provided by PRESS, we developed a new energy-saving technique called READ (Reliability and Energy Aware Distribution). Experimental results show that READ outperforms MAID [5] and PDC [26], two well-known disk array energy conservation schemes, in performance and reliability while reaching a similar level of energy-saving.

In the next section we discuss the related work and motivation. In Section 3, we describe the design of the PRESS model. The READ strategy is presented in Section 4. In Section 5, we evaluate the performance of READ. Section 6 concludes the paper with summary and future directions.

2. Related work and motivation

2.1. Energy conservation techniques

Typical energy conservation techniques for parallel disk arrays can be categorized into five broad categories as shown in Table 1.

Simply shutting down disks after a period of idle time to save energy is not feasible for parallel disk storage systems as they are normally used to serve server class workloads, where idle time slots are usually too small to justify the overhead caused by frequent spin up and spin down [4,13,42]. Therefore, power management mechanisms based on multi-speed disks like DRPM [13], Multi-speed [4], and Hibernator [42] have been proposed so that one can dynamically modulate disk speed to control the energy consumption. While two software-oriented proactive approaches, CMTPM and CMDRPM [36], can adjust the speed of disks prior to real accesses based on the access patterns provided by compiler, Hibernator [42] mainly employs disks that can spin at different speeds. The idea behind techniques in this category was inspired by technologies widely used in laptop environments [33] where multiple disk power modes are presented. Essentially, these techniques completely depend on the availability of multi-speed disks offered by disk manufacturers. Although real multi-speed (more than 2 speeds) DRPM disks are not widely available in the market yet [24], a few simple variations of DRPM disks, such as a two-speed Hitachi Deskstar 7 K 400 hard drive, have recently been produced [17].

Workload skew based energy conservation techniques concentrate the majority of workload onto a subset of a disk array so that other disks can have chances to operate in low-power modes to save energy. MAID [5] and PDC [26] are two representative energysaving schemes in this category. When conventional single-speed disks are used, both PDC and MAID can only conserve energy when the load on the server is extremely low, which is not typical for real network server workloads. However, when multi-speed disks are employed, they can significantly save energy with only a small degradation in user response time [26]. The limitation of this type of techniques is that skewed workload adversely affects disk reliability due to the load concentration. Although the authors of [26] indicate that the investigation on the reliability implications of PDC is one of their future works, no further study results on this issue have been reported in the literature to date.

Another possible solution to save energy for disk arrays is to configure the parameters of RAID. RAID Tuning [14] can achieve a good trade-off between energy and performance by carefully tuning parameters such as RAID type, number of disks, and stripe size. Similarly, PARAID [37] employs a skewed stripe pattern to adapt to the system workload by varying the number of powered disks. The advantage of these schemes is that they have high data reliability as they are built on top of RAID. The disadvantage is also obvious because they have to bear the burden of a replication based reliability mechanism to match different RAID levels [37].

Since caching is a commonly used technique for bridging the performance gap between CPU and memory so that memory access time can be significantly reduced, it is natural for one to utilize the cache to facilitate energy-saving for parallel disk systems by decreasing the number of disk accesses. While PA-LRU and PB-LRU [43] are two power-aware storage cache management policies that are LRU based, CBSM [3] periodically resizes the cache memory to control disk accesses.

Originally, disk layout or data placement strategies, which allocate all the data onto a disk array before they are accessed or dynamically migrate data between disks, were exploited to improve the overall performance in terms of mean response time of a parallel I/O system [22]. PDDL [34] and PF+ [35] extended disk layout mechanisms to make them energy-aware such that energy savings can be achieved by optimizing data placement on disks. However, they are only dedicated for array-based scientific applications. Further, to apply these algorithms, one has to modify the compiler to make it aware of the disk layout information. Finally, to better exploit existing energy-saving capabilities, their disk layout algorithms need to be combined with application code restructuring to increase the length of idle periods. This strategy demands modifications of the application's code, and thus incurs extra overhead for users.

Some of the strategies listed in Table 1 actually use a combination of several different techniques. Hibernator [42], for example, is a multi-dimensional solution that builds on a multi-speed disk mechanism (Power Management) and automatic data migration policy (Data Placement). Similarly, when utilizing multi-speed disks, MAID [5] and PDC [26] become hybrid techniques, which integrate the disk power management mechanism into workload skew technology.

2.2. Disk failure analysis and reliability estimation

Understanding disk failures and estimating disk drive reliability are two challenging tasks due to a number of reasons. First of all, a wide spectrum of disk failure-causing stresses such as age, temperature, altitude, duty cycles, activity level, and spindle start/stop frequency [10] influence disk reliability in one way or another. In addition, it is not uncommon that the root cause of a disk failure cannot be correctly identified [30,38]. Furthermore, some of the stresses might interact with each other, and thus, form a combined effect on disk reliability. For example, a higher duty cycle normally leads to a higher average disk operating temperature. Besides, a failure-causing stress may not keep a constant effect on disk reliability during disk life time. Rather, its impact varies with the time. Pinheiro et al. [27] found that the trend for higher failure rates with higher temperature is much more pronounced for old drives than for new drives. Second, disk drive reliability is typically represented by failure rate, which is highly correlated with drive models, manufacturers, and vintages [27]. Therefore, measuring the reliability of a large-scale storage system having a huge population of disks in different models becomes more complicated. Finally, the basic failure rate of a hard disk drive is a function of time and vintage rather than a constant MTBF value. Generally, it follows a stair-step distribution [19]. As such, a comprehensive investigation on large populations of disks is the only avenue to collect sufficient failure statistics in order to accurately estimate disk reliability [27]. However, large disk failure data sets from large-scale production systems are normally not publicly available [30], which prevents researchers from achieving a deep understanding of disk failures.

Existing research studies in disk failure analysis and reliability estimation can be generally categorized into two camps: vendor technical papers and user empirical reports [27]. Disk manufacturers' investigations on disk failure and reliability estimation mainly employ two technologies, mathematical modeling and laboratory testing. Cole estimated the reliability of drives in desktop computers and consumer electronics by using Seagate laboratory test data and Weibull parameters [6]. Yang and Sun introduced how Quantum made reliability predictions based on accelerated life tests and field tracking data [40]. Shah and Elerath from Network Appliance performed a series of reliability analyses based on field failure data of various drive models from different drive manufacturers [10,32]. The biggest problem for manufacturer papers is that MTBF hours are often overestimated. The cause of the problem can be attributed to the limitations of extrapolations from manufacturers' accelerated life experiments [27]. Compared with numerous vendor technical papers, there are only a very few user empirical reports so far. Schwarz et al. found a 2% disk failure rate from a total of 2489 disks deployed at the Internet Archive in the complete 2005 year based on their Archive Observatory [31]. Very recently, two pioneer investigations from Google [27] and CMU [30] opened up new perspectives for gaining a better understanding of disk failures in large-scale production systems. Schroeder and Gibson in [30] analyzed disk replacement data from several large deployments and observed a largely overstated datasheet MTBF specified by manufacturers. They found that the annual disk replacement rates in the field are usually in the range from 2% to 4%, which is much higher than manufactures' datasheet annual failure rate (e.g., 0.88% for disks with 1,000,000 h MTBF). Pinheiro et al. focused on finding how various factors such as temperature and activity level can affect disk drive lifetime [27]. Interestingly, they found a weak correlation between failure rate and temperature or activity levels, which is against the results from many previous works. Both [27,30] drew their conclusions based on field data collected from over 100,000 drives with some for an entire lifetime of five years. We believe that user empirical reports from field data are more practicable than manufacturer technical papers, which primarily rely on mathematical models and laboratory tests.

2.3. Motivation

Saving energy and maintaining system reliability, however, could be two conflicting goals. The side-effects of major energy-saving schemes on disk reliability may not be tolerated in reliability-critical applications like a mobile data center [23], where no or very small number of spare disks can be carried due to the limited space and data loss is prohibited. Thus, a better understanding of the impacts of existing disk array energy-saving schemes on disk reliability is essential for two reasons. First, only after gaining a thorough understanding on the energy-reliability conflict can one accurately assess existing energy-saving schemes

in terms of their effects on disk reliability, and then, choose an appropriate energy conservation approach for a particular application environment. Second, the complete understanding of the energy–reliability puzzle can effectively direct the design of new energy-saving techniques. Unfortunately, to the best of our knowledge, little investigation has been concentrated on this particular problem. Motivated by the importance of this largely ignored problem, in this work we study the effects of energy-saving schemes on disk reliability.

3. The PRESS model

3.1. Overview

Developing a comprehensive disk reliability model that takes all reliability-affecting factors such as age, vintage, and altitude into account is out of the scope of this work. The goal of this research is to understand the impacts of disk array energy conservation schemes on disk reliability. As such, the PRESS model embraces only major reliability-affecting factors that are influenced by disk array energy conservation schemes such as MAID [5]. All other factors like age, model, and altitude are categorized to the non-ESRRA group, and thus are ignored by PRESS. For example, although age does have influence on disk reliability [19], its reliability impact is not affected by energysaving schemes. Besides, when we compare the reliability sideeffects of various energy-saving algorithms we assume that all algorithms use the same disk array, in which all disks are identical in age. Hence, the reliability impacts incurred by age can be safely omitted in the context of this research. The same applies to other non-ESRRA factors such as altitude and vintage.

As stated in the IDEMA Standards [19], disk drive reliability is impacted each time a spindle is powered up or powered down. Similarly, when multi-speed disks are employed, power management based energy conservation schemes frequently spin up or spin down disks in order to save energy, which causes spindle motors to fail prematurely. In addition, workload skew based energy-saving strategies like PDC [26] and MAID [5], when combined with multi-speed disks, normally result in a subset of disks in a disk array always operating in higher speeds and higher activity levels due to a heavier load. Higher speeds imply higher operating temperatures, which usually lead to a high disk failure rate [19]. Thus, operating temperature, disk utilization, and disk speed transition frequency are identified as major ESRRA factors. We assume that all disks under investigation are the same in all non-ESRRA factors. Furthermore, we assume that all disks are older than 1 year, and hence, the infant mortality phenomena will not be considered in this study.

Fig. 1 depicts the overall architecture of the PRESS model. Energy-saving schemes such as DRPM [13], PDC [26], and MAID [5] inherently affect either part of the three ESRRA factors or all of them. Each of the three ESRRA factors is then fed into a corresponding reliability estimation function within the PRESS model. The PRESS model is composed of a reliability integrator module and three functions: temperature-reliability function, utilization-reliability function, and frequency-reliability function. While the former two functions were derived directly from Google's field data [27], the last one was deduced from the spindle start/stop failure rate adder suggested by the IDEMA Standards [19] and the modified Coffin–Manson model [11]. Each of the three reliability functions individually outputs its estimated reliability values in AFR (Annualized Failure Rate), which then become the inputs of the reliability integrator. The reliability integrator module translates the outputs of the three functions into a single reliability value for a disk array.



Fig. 1. Overall architecture of the PRESS model.

3.2. Operating temperature

Operating temperature has long been believed as one of the most significant factors that affect disk reliability [1,6,16,19]. High temperature was discovered as a major culprit for a number of disk reliability problems [16]. One such problem is off-track writes, which could corrupt data on adjacent cylinders. Even worse, the spindle motor and voice coil motor running at high temperatures can lead to a head crash [16]. Results from Seagate based on mathematical modeling and laboratory testing indicate that disk failure rate doubles when temperature increases by 15 °C [6]. More recently, research outcomes from Google using field data confirmed that disk operating temperature generally has observable effects on disk reliability, especially for older disks in high temperature ranges [27].

There are two different avenues to establishing a temperature-reliability relationship function. One is using mathematical modeling and laboratory testing techniques [6] and the other is employing user field data [27]. Although the two ways can provide us with temperature-reliability relationship functions with a similar trend, i.e., higher temperatures usually result in higher AFR, we selected the latter because it is a more realistic, though not perfect, way to estimate disk reliability due to sufficient amount of failure statistics from real disk deployments.

Gurumurthi et al. [13] proposed a dynamic multi-speed disk model, which can dynamically change disk speed while spinning. Consequently, a multi-speed disk could serve requests at low speeds when workload is light. Sony manufactured disk drives that are designed to operate at a small set of different rotational speeds [39]. Current commercial versions of such disk drives only support two speeds [39]. Thus, in this study, we only consider a simple type of multi-speed disks, namely, two-speed disks. We assume that the low speed mode is 3600 RPM (revolutions per minute) and the high speed mode is 10.000 RPM. It is understood that operating temperature of a disk is affected by workload characteristics and several disk drive parameters like drive geometry, number of platters, RPM, and materials used for building the drive [20]. The change of RPM, however, becomes a primary influence on a disk's temperature when all other factors mentioned above remain the same. This is because disk heat dissipation is proportional to nearly the cubic power of RPM [20]. Therefore, the increase of RPM results in excessive heat, which in turn leads to a higher temperature. Since there is no explicit information about the relationship between RPM and disk temperature, we derive temperatures of two-speed disks at 3600 and 10,000 RPM based on reported related work. Based on the IDEMA Standards [18], the temperature of the air used to cool the drives cannot result in a temperature of less than 35 °C. Besides, the experimental report in [13] indicates that on average the temperature of a hard disk drive with 5400 RPM is 37.5 °C. Therefore, our assumption that the low speed mode 3600 RPM is associated with a temperature range $[35-40 \ ^{\circ}C]$ is feasible. Also, the experimental results in [14] show that a Seagate Cheetah disk drive reaches a steady state of 55.22 $^{\circ}C$ when running at 15,000 RPM after 48 min. Considering that 10,000 RPM is only 2/3 of the disk's rotation speed, we argue that [45-50 $^{\circ}C$] is a reasonable temperature range for the high speed mode.

Now we explain why we adopted the 3-year temperature-AFR statistics from [27] as our temperature-reliability function. One can easily make the following two observations from Fig. 2(a), which is Figure 5 in [27]. First, higher temperatures are not associated with higher failure rates when disks are less than 3 years old. Second, the temperature effects on disk failure rates are salient for the 3-year-old and the 4-year-old disks, especially when temperatures are higher than 35 °C. The authors of [27] explain the reason of the first observation is that other effects may affect failure rates much more strongly than temperatures do when disks are still young. However, we have a different interpretation of this phenomenon. We argue that higher temperatures still have strong negative effects on younger disks as they do on older disks. The impacts of higher temperature on younger disks do not immediately turn out to be explicit disk failures just because the impacts-to-failure procedure is essentially an accumulation process and it takes some time. After all, higher temperatures make electronic and mechanical components of disks more prone to fail prematurely [16]. The second observation, i.e., obvious higher failure rates associated with higher temperature ranges for 3-year-old disks, supports our explanation because earlier high temperature impacts on disks are eventually transformed into disk failures after one or two years. Therefore, we ignore the temperature-AFR results in [27] for disks younger than 3 years as they hide the temperature impacts on disk reliability. Although both 3-year-old disks and 4-year-old disks exhibit a high correlation between higher temperatures and higher failure rates, we finally decided to select 3-year-disk temperature-AFR data as the foundation of our temperature-reliability function. The primary reason is that the relationship between higher temperatures and AFR for 3-year-old disks fully demonstrates that higher temperatures have a prominent influence on disk failure rates because after 2-year higher temperature "torture" an observable number of disks fail in the third year. Apparently, these disk failures, which are originated in the first two years, should be included in the third year's AFR. On the other hand, the 4-yearold disk results substantially lose the "hidden" disk failures, and therefore are not complete. Our temperature-reliability function (see Fig. 2(b)) is based on the results of 3-year old disks in [27] within the temperature range [15–70 °C]. It shows that when disks are running in the low speed mode within the temperature range [35–40 °C], the AFR is about 6.5%. If the disks are operating in high speed mode falling in the temperature scope [45–50 °C], the corresponding AFR is around 15%.



Fig. 2. (a) Temperature impacts on AFR in [27]. (b) The temperature-reliability function in the PRESS model.



Fig. 3. (a) Utilization impacts on AFR in [27]. (b) The utilization-reliability function in the PRESS model.

3.3. Disk utilization

Disk utilization is defined as the fraction of active time of a drive out of its total power-on-time. Since there is not enough detail in their measurements, the researchers of [27] measured utilization in terms of weekly averages of read/write bandwidth for each drive and roughly divided them into three categories: low, medium, and high. Still, they found that using the number of I/O operations and bytes transferred as utilization metrics provided very similar results [27]. Thus, we conclude that it is feasible to take the average bandwidth metric as the utilization metric because the number of I/O operations and bytes transferred of a disk are proportional to disk active time. Therefore, in our utilization-reliability function we use the utilization metric in the range [25%–100%] instead of low, medium, and high employed in Figure 3 of [27]. We define low utilizations as utilizations in the range [25%-50%]. Similarly, a medium utilization is defined as a utilization within the scope [50%-75%], whereas a high utilization falls in the range [75%-100%].

The relationship between utilization and disk reliability has been investigated previously [1,6,27,40]. A conclusion that higher utilizations in most cases affect disk reliability negatively has been generally confirmed by two widely recognized studies. One is a classical work from Seagate, which utilized laboratory testing and mathematical modeling techniques [6]. The other is a new breakthrough, which analyzes the utilization impacts on disk reliability based on field data from Google [27]. The authors of [27] measured 7 age groups of disks (3-month, 6-month, 1-year, 2-year, 3-year, 4-year, 5-year, see Figure 3 in [27]) and found that only the 3-year-old group exhibits an unexpected result, i.e., low utilizations result in a slightly higher AFR than higher utilizations do. The two explanations for this "bizarre" behavior provided by [27] are not convincing in our view. Their first explanation is the survival of the fittest theory. They speculate that the drives that survive the infant mortality phase are the least susceptible to the failures caused by higher utilizations, and result in a population that is more robust with respect to variations in utilization levels [27]. If this is the case, they cannot explain

why the results from the 4-year-old disk group and the 5-yearold disk group immediately restore the "wired" behavior to a "normal" one, i.e., higher utilizations correlate to higher AFR. The second explanation they made is that previous results such as [6] can only better model early life failure characteristics, and thus, it is possible that longer term population studies could discover a less significant effect later in a disk's lifetime. Again, if this is true, they cannot explain why we still see a noticeable higher utilization with higher AFR behavior for disks in their ages 4 and 5. In fact, their second explanation conflicts with their observation that only very young and very old age groups show the expected behavior. Based on our observations in Fig. 3(a), we argue that a reasonable explanation for this unexpected behavior is that disk drives in their middle ages (2 or 3 years) are strong enough in both electronic and mechanical parts to resist the effects of higher utilizations. Therefore, AFR of disks in these two age groups has little correlation with utilization. Our speculation is supported by the evidence that failure rates of different utilization levels are very close to each other for disks in these two age groups and failure rate distribution exhibits some randomness. We selected the results from the 4-year-old disk group as our utilization-reliability function mainly because (1) we only consider disks older than 1 year; (2) results from 2-year and 3-year groups cannot provide any explicit utilization impacts on disk reliability although much previous research confirms that these impacts do exist; (3) 5year results are less useful because disks normally only have five year warranty; and (4) the results from 4-year disks match the reliability versus duty cycle outcomes of [6]. Based on the results from the 4-year disk group, we built our utilization-reliability function as shown in Fig. 3(b). While utilizations in the range [25%–50%] lead to an AFR around 2.5%, utilizations between 50% and 75% result in an AFR of about 3.5%. The high utilizations in the range [75%-100%] cause an AFR of 3.8%.

3.4. Disk speed transition frequency

The disk speed transition frequency (hereafter called frequency) is defined as the number of disk speed transitions in one



Fig. 4. (a) Start/Stop failure rate adder [19]. (b) The frequency-reliability function in the PRESS model.

day. Establishing a frequency–reliability function is the most difficult task in this research primarily because multi-speed disks have not been largely manufactured and deployed. Thus, no result about the impacts of frequency on disk reliability has been reported so far. Although the applications of multi-speed disks are still in their infancy, we believe that they will no doubt have a huge impact on energy-saving for disk-based storage systems in the not-so-distant future. Therefore, now it is the time to obtain a basic understanding of the relationship between frequency and reliability. Our frequency–reliability function is built on a combination of the spindle start/stop failure rate adder suggested by IDEMA [19] and the modified Coffin–Manson model.

We start our investigation on this challenging issue from a relevant disk usage pattern parameter, namely, spindle start/stop rate (SSSR), defined as the times of spindle start/stop per month [10,19]. The rationale behind this is that disk speed transitions and spindle start/stops essentially generate the same type of disk failure mode, spindle motor failure, though with different extents. A disk reliability report discovered that each spindle start-and-stop event causes some amount of fatigue to occur at the heads and the spindle motor [29]. In fact, spindle motor failure is one of the most common disk drive failure modes [21]. That is why disk drive manufacturers normally set 50,000 as the start/stop cycle limit and suggest no more than 25 power cycles per day to guarantee specified performance. A disk speed transition event could cause a similar reliability issue as a spindle start/stop occurrence does because speed transitions incur some amount of fatigue, noise, heat dissipation and vibration as well [21]. Among these side-effects, fatigue is a dominant disk reliability-affecting factor [21]. Therefore, we will focus on disk reliability impact brought by fatigue while ignoring all other factors like vibration. We believe, however, the degree of reliability impacts caused by speed transitions is relatively lower than that caused by spindle start/stops. The reason is two-fold. First, during a start up process, a spindle has to increase its speed from zero to maximum. However, a speed transition event, e.g., from a low speed to a high speed, only needs to promote the spindle's speed from its current value to an immediate higher value. Therefore, the costs of a speed transition between two contiguous speed levels in terms of energy consumption and time are less than that of a spindle start/stop, which in turn brings a disk drive less heat dissipation, a main reason for fatigue. Second, there is no salient peak power issue associated with speed transitions. It is understood that peak power within a short period of time is detrimental to disk reliability [21].

Both start/stop events and disk speed transitions incur temperature cycling, the main cause of fatigue failures [9]. The damage caused by temperature cycling accumulates each time a hard disk drive undergoes a power cycle or a speed transition. Such cycles induce a cyclical stress, which weakens materials and eventually makes the disk fail [9]. We utilize the modified Coffin–Manson model (Eq. (1)) because it is a widely-used model, which works very well for failures caused by material fatigues due to cyclical stress [11]. It evaluates the reliability effects of cycles of stress or frequency of change in temperatures. The Arrhenius equation involved describes the relationship between failure rate and temperature for electronic components (Eq. (2)).

The spindle start/stop failure rate adder curve presented by IDEMA is re-plotted as Fig. 4(a). It indicates, for example, a start/stop rate of 10 per day would add 0.15 to the AFR for disks older than one year. Since IDEMA only gives the curve in a start/stop frequency range [0-350] per month, we extend it to [0-1600] per day using quadratic curve fitting technique. We derive our frequency-function based on Fig. 4(a) and the modified Coffin–Manson model, which is given by Eq. (1) below:

$$N_f = A_0 f^{-\alpha} \Delta T^{-\beta} G(T_{\max}), \tag{1}$$

where N_f is the number of cycles to failure, A_0 is a material constant, f is the cycling frequency, ΔT is the temperature range during a cycle, and $G(T_{\text{max}})$ is an Arrhenius term evaluated at the maximum temperature reached in each cycle. Typical values for the cycling frequency exponent α and the temperature range exponent β are around -1/3 and 2, respectively [11]. The term $G(T_{\text{max}})$ can be calculated using the following Arrhenius equation [11]:

$$G(T) = A e^{(-E_a/KT)},$$
(2)

where A is a constant scaling factor, E_a is the activation energy, K is the Boltzmann's constant (i.e., 8.617×10^{-5}), and T is the temperature measured in degrees Kelvin (i.e., 273.16 + degrees in Celsius) at the point when the failure process takes place.

We first demonstrate how we derive the value of $G(T_{max})$ using Eq. (2). Since the maximum disk operating temperature is set to 50 °C when a disk is running at its high speed, T_{max} is equal to 273.16 + 50 = 323.16 K. Also, E_a is suggested to be 1.25 [11]. Therefore, $G(T_{\text{max}}) = A * 3.2275 \times 10^{-20}$. Since the suggested daily power cycle limit is 25, we set f equal to 25. Also, the temperature gap from an ambient temperature 28 °C to the maximum operating temperature 50 °C is 22 °C, which means that ΔT is equal to 22. Besides, we know that the maximum number of power cycles specified in a disk datasheet is normally 50,000. We let N_f be 50,000. Consequently, based on Eq. (1), we obtain $A * A_0 = 2.564317 \times 10^{26}$. Now we calculate N'_f , the number of speed transitions to failure assuming that the number of speed transitions per day is 25. Here, the temperature T_{max} is set to 45 °C, the midway value of the low temperature 40 °C and the high temperature 50 °C (see Section 3.2). The reason is that speed transition is bi-directional in the sense that a speed transition could either increase or decrease disk temperature. Now ΔT in Eq. (1) is equal to 10 because this is the gap between the low temperature range and the high temperature range (see Section 3.2). Based on Eq. (1) and the calculated value of $A * A_0$, we conclude that N'_f is equal to 118529, the number of disk speed transitions to failure. We view this as strong evidence that a disk speed transition can cause about 50% of the effect on reliability incurred by a spindle start/stop. Therefore, we scale down the spindle start/stop failure



Fig. 5. (a) The PRESS model at 40 °C temperature. (b) The PRESS model at 50 °C temperature.

rate adder curve (Fig. 4(a)) by half and change the unit of the X axis to times per day to obtain our frequency–reliability function shown in Fig. 4(b). The mathematical expression based on quadratic curve fitting for the reliability-frequency function is Eq. (3), where R is the reliability in AFR and f is the disk speed transition frequency.

 $R(f) = 1.51e^{-5}f^2 - 1.09e^{-4}f + 1.39e^{-4}, f \in [0, 1600].$ (3)

3.5. PRESS it all together

Since we 3-dimensional people have no 4-dimensional perspective, we present two 3-dimensional figures to represent the PRESS model at operating temperature 40 °C (Fig. 5(a)) and 50 °C (Fig. 5(b)), respectively. In Section 3.2, we justified why it is reasonable to set the temperature range [35–40 °C] for disk speed 3600 RPM and the temperature range [45–50 °C] for disk speed 10,000 RPM. Within these feasible temperature range settings, we suppose that disks at nlow speed have operating temperature 40 °C, whereas disks at high speed are at 50 °C. After obtaining the AFR for each disk in a disk array, the reliability integrator module outputs the AFR of the least reliable disk as the overall reliability for the entire disk array. This is because the reliability level of a disk array is only as high as the lowest reliability level of a single disk in the array.

The PRESS model yields several important insights on how to make trade-offs between energy-saving and reliability when developing energy conservation techniques for disk array systems. First, disk speed transition frequency is the most significant reliability-affecting factor among the three ESRRA factors. Based on our estimation in Section 3.4, the number of disk speed transitions should be limited to less than 65 (118529/5/365 \approx 65) per day in order to guarantee a 5-year performance warranty. Thus, it is not wise to aggressively switch disk speed to save some amount of energy. We argue that the high AFR caused by a high speed transition frequency would cost much more than the energysaving gained. Normally, the value of lost data plus the price of failed disks substantially outweigh the energy-saving gained. Thus, it is not worthwhile for disk arrays to save energy by frequently switching disk speed. This argument has been validated through our experimental results presented in Section 5. For example, on average MAID can save 2.3% energy consumption compared with READ in a heavy workload condition (Fig. 9(b)) when the number of disks changes from 12 to 16. However, the average AFR of MAID is 25.2% higher than that of READ in the same condition (Fig. 7(b)). Next, operating temperature is the second most significant reliability-affecting factor. A high temperature can be caused by long time running at high speed. Hence, workload skew based energy-saving schemes need to rotate the role of workhorse disks regularly so that the scenario that a particular subset of disks is always running at high temperature can be avoided. Finally, since the AFR differences between high utilizations and medium utilizations are slim, an uneven utilization distribution should not be overly concerning.

3.6. Validation of PRESS

In this section, we conduct a preliminary validation study on the PRESS model. Since there is no report about the impacts of frequency on disk reliability in the literature, we only validate the temperature–reliability function and the utilization–reliability function here. Still, we believe that the frequency–reliability function is useful because it is reasonably derived from the modified Coffin–Manson model, which has been successfully used to model materials' fatigue failures due to repeated temperature cycling as a device is turned on and off [9,11]. We will validate the frequency-reliability function once field data from a large-scale deployment of multi-speed disk drives are available.

We validate the temperature-reliability function and the utilization-reliability function by comparing them with the findings of other studies [1,6,19,30]. We found that although the two functions are determined from empirical observations [27], they are generally consistent with the results from other field data analysis [19,30] and standard laboratory tests [1,6]. We take this consistency as strong evidence that the two functions are valid and reasonably accurate.

A sample engineering datasheet in the IDEMA Standards [19] suggests that the development estimated average base failure rate in %/1000 h for disks older than 1 year is 0.15. In other words, the estimated average base AFR is 0.15 * 24 * 365/1000 =1.314%. Considering the impacts of temperature on AFR in terms of AFR multiplier, the temperature range [35-40 °C] on average brings the AFR multiplier 1.4 and the temperature range [45–50 °C] on average causes the AFR multiplier 2.6 according to the temperature-environment factor figure provided in [19]. As a result, the estimated average datasheet AFR values for disks in the temperature range [35-40 °C] and [45-50 °C] are 1.314% * 1.4 = 1.84% and 1.314% * 2.6 = 3.42%, respectively. Schroeder and Gibson discovered that the annual disk replacement rates in the field are usually in the range from 2% to 4% for disks with manufactures' datasheet annual failure rate 0.88% (i.e., 1000,000 h MTBF) [30]. The implication of their findings is that an estimated datasheet AFR needs to be on average amplified by around 3.4 times [(2% + 4%)/2/0.88% = 3.4]. Thus, the expected field average AFR values for the temperature ranges [35–40 °C] and [45–50 °C] are 1.84%*3.4 = 6.3% and 3.42%*3.4 = 11.6%, respectively. These results are consistent with the temperature-reliability function shown in Fig. 2(b) (see Section 3.2).

Duty cycle is defined as the fraction of time a drive is active out of the total powered-on time, which is also called utilization by the literature [27]. Based on laboratory tests in Seagate [1,6], Cole found that the average AFR multiplier factors for duty cycle ranges [25%–50%], [50%–75%], and [75%–100%] are around 0.6, 0.8, and 0.9, respectively (see Figure 10 in [1]). Thus, the estimated average AFR values for these three duty cycle ranges are 1.314% * 0.6 = 0.78%, 1.314% * 0.8 = 1.05%, and 1.314% * 0.9 = 1.18%, respectively. Consequently, the expected field average AFR values

for the three duty cycle ranges are 0.78% * 3.4 = 2.65%, 1.05% * 3.4 = 3.57%, and 1.18% * 3.4 = 4.01%, respectively. The three expected field AFR values are very close to the AFR values (2.5% for utilization range [25%–50%], 3.5% for utilization range [50%–75%], and 3.8% for utilization range [75%–100%]) plotted in Fig. 3(b) (see Section 3.3).

4. The READ strategy

With the light shed by the PRESS model on design of new reliability-aware energy-saving techniques for disk arrays, in this section we first present a general idea of our READ approach, which is followed by a detailed algorithm description, as well as a complexity analysis of READ.

4.1. READ reliability and energy aware distribution

Several previous studies [7,12] show that the distribution of web page requests generally follows a Zipf distribution [22] where the relative probability of a request for the *i*'th most popular page is proportional to $1/i^{\alpha}$, with α typically varying between 0 and 1. Further, they discover that the request frequency and the file size are inversely correlated [7,12], i.e., the most popular files are typically small in size, while the large files are relatively unpopular. Inspired by the observations of this highly skewed data popularity distribution, two traditional energy-saving techniques, MAID [5] and PDC [26], concentrate the majority of workload onto a subset of a disk array so that other disks can have chances to operate in lowpower modes to save energy. PDC dynamically migrates popular data to a subset of the disks so that the load becomes skewed towards a few of the disks and others can be sent to low-power modes [26]. Since only a small portion of data would be accessed at a given time, the idea of Massive Array of Idle Disks (MAID) [5] is to copy the required data to a set of "cache disks" and put all the other disks in low-power mode. Later accesses to the data may then hit the data on the cache disk(s). A common goal of both PDC and MAID is to increase idle times by rearranging data among the disk array and lower the disks' speed down [26]. Neithere of the two algorithms applied any mechanisms to limit the reliability impacts introduced by them.

Our READ strategy is motivated by data popularity locality as well and it employs data redistribution and multi-speed disks. We adopted several similar assumptions that PDC used. We also assume that each request accesses an entire file, which is a typical scenario for Web, proxy, ftp, and email server workloads [26]. In addition, the distribution of requests generally follows a Zipflike distribution with α in the range [0, 1]. Also, each file is permanently stored on one disk and neither striping nor mirroring is used [26]. We decide not to use striping for two reasons. One is that we want to make the comparisons between READ and the two conventional algorithms in a fair manner as they did not employ striping. The other is that the average file sizes in the real web workload are much smaller than a normal striping block size 512 KB. Further, no requests can be served when a disk is switching its speed. The general idea of READ is to control disk speed transition frequency based on the statistics of the workload so that disk array reliability can be guaranteed. Also, READ employs a dynamic file redistribution scheme to periodically redistribute files across a disk array in an even manner to generate a more uniform disk utilization distribution. A low disk speed transition frequency and an even distribution of disk utilizations imply a lower AFR based on the PRESS model.

4.2. Description of the algorithm

The set of files is represented as $F = \{f_1, \ldots, f_u, f_v, \ldots, f_m\}$. A file f_i ($f_i \in F$) is modeled as a set of rational parameters, e.g., $f_i = (s_i, \lambda_i)$, where s_i, λ_i are the file's size in Mbyte and its access rate. In the original round of file distribution, READ orders the files in terms of file size because we assume that the popularity in terms of access rate of a file is inversely correlated to its size. And then READ splits the file set into two subsets: popular file set F_p = $\{f_1, \ldots, f_h, \ldots, f_u\}$ and unpopular file set $F_u = \{f_v, \ldots, f_c, \ldots, f_m\}$ $\{F = F_p \cup F_u \text{ and } F_p \cap F_u = \emptyset\}$. Next, a disk array storage system consists of a linked group $D = \{d_1, \ldots, d_e, d_f, \ldots, d_n\}$ of *n* independent 2-speed disk drives, which can be divided into a hot disk zone $D_h = \{d_1, \ldots, d_h, \ldots, d_e\}$ and a cold disk zone $D_c = \{d_f, \ldots, d_c, \ldots, d_n\}$ $(D = D_h \cup D_c \text{ and } D_h \cap D_c = \emptyset).$ Disks in the hot zone are all configured to their high speed modes, which always run in the high transfer rate t^h (MB/s) with the high active energy consumption rate p^h (J/MB) and the high idle energy consumption rate i^h (J/s). Similarly, disks in the cold zone are set to their low speed modes, which continuously operate in the low transfer rate t^{1} (MB/s) with the low active energy consumption rate p^{l} (J/MB) and the low idle energy consumption rate i^{l} (J/s). All disks have the same capacity c.

READ places popular files onto the hot disk zone and unpopular files onto the cold disk zone. The ratio between hot disk number and cold disk number in a disk array is decided by the load percentages of popular files and unpopular files in the whole file set. The load of a file f_i is defined as $h_i = \lambda_i \cdot sv_i$, where sv_i, λ_i are the file's service time and its access rate. Since we assume that each request sequentially scans a file from the beginning to the end, sv_i is proportional to s_i , the size of file f_i . Thus, the load of file f_i can also be expressed as $h_i = \lambda_i \cdot s_i$. Besides, we assume that the distribution of file access requests is a Zipf-like distribution with a skew parameter $\theta = \log \frac{A}{100} / \log \frac{B}{100}$, where *A* percent of all accesses are directed to *B* percent of file [22]. The number of popular files in *F* is defined as $|F_p| = (1 - \theta) * m$, where *m* is the total number of files in *F*. Similarly, the number of unpopular files and the number of unpopular files in *F* is defined as δ

$$\delta = (1 - \theta)/\theta. \tag{4}$$

The ratio between the number of hot disks and the number of cold disks is defined as γ , which is decided by the ratio between the total load of popular files and the total load of unpopular files:

$$\gamma = \sum_{i=1,f_i \in F_p}^{(1-\theta)*m} h_i \bigg/ \sum_{j=1,f_j \in F_u}^{\theta*m} h_j.$$
⁽⁵⁾

Fig. 6 depicts the READ algorithm. READ assigns sorted popular files in F_p onto the hot disk zone in a round-robin manner with the first file (supposed the most popular one) onto the first disk, the second file onto the second disk, and so on. Similar file assignment strategy is applied for sorted unpopular files in F_u onto the code disk zone. After all files in F have been allocated, READ launches an Access Tracking Manager (ATM) process, which records each file's popularity in terms of number of accesses within one epoch in a table called File Popularity Table (FPT).

The FPT table with the latest popularity information for each file will be used later by the File Redistribution Daemon (FRD). At the end of each epoch, FRD re-orders all files based on their access times recorded during the current epoch in the FPT table and then redefines popular file set F_p and unpopular file set F_u accordingly. A hot file will be migrated to the cold disk zone if its new position in the entire re-sorted file set is out of the newly defined hot file range. It will stay in the hot zone, otherwise. Similarly, a previous cold file will be migrated to the hot disk zone if its new ranking is within the new hot file scope.

4.3. Time complexity of READ

Before qualitatively comparing our scheme with the two existing algorithms, we demonstrate the worst-case time complexity of the READ algorithm in Theorem 1.

Input: A disk array D with n 2-speed disks, a collection of m files in F, an epoch P, idleness threshold H, a disk maximum allowed provide transitions per day S speed transitions for each disk $T(n)$ the skew perspects θ
speed transitions per day s, speed transition times to each disk $f(n)$, the skew parameter σ
1 Lies Eq. 4 to compute the number of nonular files and the number of unpopular files
2. Lise Eq. 5 to compute the harmost popular mes and the number of any popular mes
3. Hot disk number $p_{n,n}$ and disk number $Q_{n,n}$ and $Q_{n,n}$ and $Q_{n,n}$ and $Q_{n,n}$ and $Q_{n,n}$
$HD = \frac{1}{2} + \frac{1}{2}$, and allow the most observed in HD , a_{h} is a constant of the second
4. Configure HD of <i>n</i> disks to high speed mode and set CD of <i>n</i> disks to low speed mode
5. Sort all files in file size in a non-decreasing order
6 Assian all popular files onto the bot disk zone
7. Assign all unpopular files onto the cold disk zone
8. for each epoch P do
9. Keep tracking number of accesses for each file
10. Re-sort files in number of accesses in current epoch
11. Re-calculate θ and re-tag popular and unpopular
12. for each previously hot file that becomes cold do
13. Migrate it to the cold disk zone
14. Update its record in the allocation scheme <i>X</i>
15. end for
16. for each previously cold file that becomes popular do
17. Migrate it to the hot disk zone
18 Update its record in the allocation scheme X
19. end for
20. for each disk $d_i \in D$ do
21. If $S/2 \le I(a_i)$ // Still has room in terms of disk speed transitions to spin down
22. H=2H; // Double the laleness threshold H to reduce future disk speed transitions
24. end for 25. and for

Fig. 6. The READ strategy.

Table 2

System parameters. Description Value Description Value Disk model Seagate Cheetah ST39205LC Standard interface SCSI Storage capacity 9.17 GB High rotational speed 10,000 RPM Number of platters Transfer rate at high speed 31 MB/s 1 High speed working temperature 50 °C High speed Idle power 5.26 W Read energy (8 K) 0.061 J Inactive power 1.86 W Spin down energy 28.25 J Read unit 8 K Spin up energy 65.91 J Spin down time 5.62 s 3.06 s Spin up time Low rotational speed 3600 RPM Low speed working temperature 40 °C Low speed idle power 2.17 W Read energy (8 K) 0.043 J Transfer rate at low speed 9.3 MB/s Maximum speed transitions per day S = 40Epoch 2 h Initial skew parameter $\theta = 0.296$ Idleness threshold 17.9 s

Theorem 1. Given a parallel disk array system $D = \{d_1, d_2, ..., d_j, ..., d_n\}$, a collection of files represented by a file set $F = (f_1, f_2, ..., f_i, ..., f_m)$, and maximum allowed disk speed transition number S for each disk, the worst-case time complexity of READ is O((2k + 2)m + (k + 1)mlgm + knS)), where m is the number of files in F, n is the number of disks in the system D, and k is the number of total epochs during the execution of READ.

Proof. It takes O(m) time to derive an appropriate value of γ based on Eq. (5) (see Step 2). Step 5 takes O(mlgm) to sort the file set *F*. The initial file assignment process (Steps 6–7) costs another O(m). Within each epoch *P*, the time spent on tracking number of accesses for each file is O(m) (Step 9). In addition, each re-sorting in Step 10 takes O(mlgm). The worst case time complexity for file redistribution is O(m) (Steps 12–19) assuming that each file needs to be migrated, which is almost impossible. Similarly, the worst case for Steps 20–24 is that each disk is spun down *S* times, which implies O(nS) worst case time complexity. Other steps simply take O(1). Further, we assume that there are a total of *k* epochs during the execution of READ. Thus, the worst-case time complexity is: O(2m + mlgm) + k(O(2m) + O(mlgm) + O(nS)) = O((2k + 2)m + (k + 1)mlgm + knS).

5. Performance evaluation

In this section, we present results of experimental simulations using real-world traces. We compare our READ strategy with two traditional disk-array energy-saving techniques PDC and MAID. We first outline the execution-driven simulator and experimental setup. Second, three widely-used real-world traces are introduced. Finally, we analyze results from trace-driven simulations.

5.1. Experimental setup

We developed an execution-driven simulator that models an array of 2-speed disks. The same strategy used in [26] to derive corresponding low speed mode disk statistics from parameters of a conventional Cheetah disk was adopted in our study. The main characteristics of the 2-speed disk and major system parameters are shown in Table 2.

The performance metrics by which we evaluate system performance include:

• Mean response time: average response time of all file access requests submitted to the simulated 2-speed parallel disk storage system.

Table 5		
Statistics of the	four real	traces.

Table 2

Trace name	Files	Requests	Ave. arrival interval (ms)	Min file size (bytes)	Max file size (bytes)	Mean file size (bytes)
Clarknet-HTTP log	10,798	186,397	463.5	25	4,571,168	13,097
World Cup 98-05-09	4,079	1,480,081	58.4	4	2,891,887	20,021
World Cup 98-06-11	4,261	59,201,342	1.46	4	3,248,697	20,086
Auspex	43,579	2,460,587	37.1	8192	402,956,720	12,552



Fig. 7. A comparison of the three algorithms in terms of reliability.

- Energy consumption: energy consumed by the system during the process of serving all requests.
- AFR: annualized failure rate of a disk array. Each disk has an AFR calculated by the PRESS model. The highest one is used to designate the AFR of the entire disk array.

We evaluate the three algorithms by running trace-driven simulations over three Web I/O traces (ClarkNet-HTTP log [4], WorldCup98-05-09, and WorldCup98-06-11 [2]) and the Auspex trace [8], which have been widely used in the literature. ClarkNet-HTTP log was collected by ClarkNet, an Internet service provider, for a week from 09/04/95 to 01/10/95 with a total of 3,328,587 requests. The frequency of file access in the trace follows a Zipflike distribution. Since the simulation times in our experiments are much shorter compared with the time span of the ClarkNet-HTTP trace, we only choose one day (09/04/95) data, which has 186,397 requests. Similarly, we select two days' data, 05/09/98 and 06/11/98, from the WorldCup98 trace, which presents one of the largest Web workloads analyzed so far [2]. While the ClarkNet-HTTP log presents our simulated disk array system a light workload, the two WorldCup98 traces introduce a heavy and an extremely heavy workload condition, respectively. While the three Web traces are read-only, the Auspex trace, which was originated from Berkeley, has both reads and writes [8]. Therefore, we can evaluate the three algorithms under various workload conditions. Table 3 shows the relevant information of the real traces.

5.2. Reliability

We conduct our performance evaluation of the three energysaving algorithms on a simulated platform of a disk array consisting of 6–16 disks. Across all given workloads and disk numbers in our experiments, the READ algorithm almost consistently outperforms MAID and PDC algorithms in reliability by up to 39.7% and 57.5%, respectively. When the lightest workload ClarkNet-HTTP log is applied onto the system, the average AFR for READ is 17.2%, whereas it is 23.5% and 36.6% for MAID and PDC, respectively (Fig. 7(a)). We attribute the substantial improvement of the READ algorithm in terms of reliability to the very limited number of disk speed transitions it incurred. When workload is very light, average request arrival interval times for unpopular files are larger than the idleness threshold. Therefore, PDC and MAID have many opportunities to spin down disks to try to save energy. On the contrary, READ constrains each disk's number of speed transitions so

that it cannot be larger than S, which was set to 40 in our study. READ accomplishes this by gradually enlarging the idleness threshold value. In our implementation, we simply double the idleness threshold value once READ finds that a disk's current number of speed transitions reaches half of S. The large number of speed transitions brings MAID and PDC high values of AFR. With the increase in workload intensity (Fig. 7(b)), the improvements of READ in AFR decrease when the disk number varies from 6 to 10. This is because a heavy workload gives MAID and PDC fewer chances to spin down, and thus the numbers of speed transitions for the two algorithms reduce. A reduced number of disk speed transitions results in a decreased value of AFR. However, when disk number increases, average load on each disk decreases, and thus PDC frequently spins down disks once again, which makes its AFR rise. When the extremely heavy workload WorldCup98-06-11 is introduced to the system, READ ties with MAID and PDC in AFR because none of them can have chances to spin down disks no matter how many disks are included in the disk array (Fig. 7(c)).

5.3. Performance

The READ algorithm delivers much shorter mean response times in all cases (Fig. 8) primarily due to its very low number of disk transitions. Interestingly, when the workload is not extremely heavy, the mean response times of MAID and PDC increase with an enlarged number of disks (Fig. 8(a) and (b)). The reason behind this is that the increased number of disks lowers the load for each disk, and thus MAID and PDC have more opportunities to spin down disks. A large number of disk spin downs implies an equivalent number of spin ups, which introduces a considerable delay in responding to requests. In an extremely heavy workload, the response times of all the three algorithms decrease with an increased number of disks (Fig. 8(c)). This observation is expected because a larger number of disks can reduce each disk's load, and thus improve mean response times. More importantly, in this situation all three algorithms almost have no chances to spin down disks even with an increased number of disks, which can noticeably boost system performance. Still, READ improves mean response times on average 77.8% and 72.2% (Fig. 8(c)).

5.4. Energy conservation

In terms of energy conservation, READ performs obviously better than the two baseline algorithms in a light workload



Fig. 8. A comparison of the three algorithms in terms of performance.



Fig. 9. A comparison of the three algorithms in terms of energy consumption.



Fig. 10. A comparison of the three algorithms in the Auspex trace.

condition (Fig. 9(a)). In the ClarkNet trace experiments, on average READ results in 4.8% and 12.6% less energy consumption compared with MAID and PDC, respectively. Obviously, all the algorithms' energy consumption goes up when the number of disks is increased as more disks consume more energy. One important observation is that a large number of disk spin downs does not necessarily bring us more energy savings. On the contrary, a disk spin down can cause more energy consumption if the idle time is not long enough to compensate for the energy cost during disk spin down and spin up. This conjecture is demonstrated by the high energy consumption of MAID and PDC in Fig. 9. We notice that our READ performs slightly worse than MAID in energy consumption in a heavy workload condition (Fig. 9(b)) when the number of disks changes from 12 to 16. The reason is that MAID still has disk spin downs when the disk number increases and these disk spin downs indeed bring energy conservation because in most cases the idle times are long enough to compensate for disk transition energy cost. On the other hand, our READ algorithm has no disk spin downs, and thus disks are always running at high speed. However, in a tremendously heavy workload condition, the energy consumption of MAID becomes the highest one among the three. This is mainly because MAID uses the hot disk zone as a "cache"

with a simple LRU replacement policy and does not keep track access pattern statistics. Hence, some requests cannot hit "cache" disks (high speed disks) and are served by the low speed disks, which consume more energy.

5.5. Read-write trace simulations

To evaluate the performance of READ in a workload scenario where both reads and writes are presented, we conducted one more group of experiments using the Auspex trace [8]. The average number of requests per second is 27 in the Auspex trace. We found that there are only around 20 speed transitions per day per disk during the simulations. Since the average request size is small (i.e., 9565 bytes/request), we cannot see a significant difference among the three algorithms in terms of reliability (see Fig. 10(a)). As the disk number increases, the requests with the high access rate can be distributed to more disks and be processed more in parallel (Fig. 10(b)). Apparently, the energy consumption increases (Fig. 10(c)). In all cases of Fig. 10, READ outperforms MAID and PDC in both performance and energy consumption. The conclusion is that READ still performs well in workload conditions where read and write requests coexist.

6. Conclusions

Numerous energy-saving techniques have been proposed to significantly reduce storage systems' energy consumption while maintaining a good performance. Unfortunately, power management based and workload skew based energy-saving schemes, two prevalent categories of energy conservation techniques for disk arrays, inherently affect disk array reliability. Hence, a thorough understanding of the relationship between energy-saving techniques and disk reliability is essential.

How to comprehensively and accurately measure reliability impacts caused by all possible factors is still an open question [10,14,23,28,34,41] and out of the scope of this research. The PRESS model is only one step towards finding a way to quantitatively approximate the reliability effects imposed by the three ESRRA factors. We believe that our model is useful due to the following two reasons. First, our temperature-reliability function and utilization-reliability function come from a state-ofthe-art work [27], which studies the impacts of the two factors on disk reliability based on field data from a large disk population over 5 years. More importantly, the two functions are generally consistent with results from previous studies [1,6,19,30]. Second, although the PRESS model is not absolutely objective and complete as it requires a number of simplified assumptions and takes only three factors into account, it can be used to capture the differences in terms of reliability among various energy-saving algorithms, which are evaluated under the same workload on the same disk array with the same environmental conditions such as altitude. After all, understanding energy-saving schemes' relative performance in reliability is more important than obtaining their absolute reliability values in the context of this work.

In this paper, we establish an empirical reliability model PRESS (Predictor of Reliability for Energy-Saving Schemes), which can be utilized to estimate reliability impacts caused by the three ESRRA factors. With the assistance of the PRESS model, system administrators can quantitatively compare existing energy-saving schemes in terms of their impacts on disk array reliability, and thus choose the most appropriate one for their applications. Besides, energy-saving technique designers can develop new energy conservation schemes, which are able to achieve a good balance between energy-saving and system reliability. Further, with the light shed by the PRESS model, we develop and evaluate a novel energy-saving strategy with reliability awareness called READ (Reliability and Energy Aware Distribution). The READ strategy exploits popularity locality of I/O workload characteristics, which is common in real workloads such as web server applications. Our trace-driven experimental results show that when workload is not extremely heavy the READ strategy results in an average 24.9% and 50.8% reliability improvement compared with MAID and PDC, respectively. Meanwhile, in terms of energy consumption, READ in most cases still outperforms the two traditional approaches. READ delivers a much better performance in mean response time.

Future directions of this research can be performed in the following directions. First, we will extend our scheme to a fully dynamic environment, where file access patterns can dramatically change in a short period of time. As a result, a high file redistribution cost may arise as the number of file migrations increases substantially. One possible solution is to use a file replication technique. Second, we intend to enable the READ scheme to cooperate with the RAID architecture, where files are usually striped across disks in order to further reduce the service time of a single request. For the web server environment, the target application domain of this work, files are usually very small (less than 50 KB), so striping is not crucial. However, for large files such as video clips, audio segments, and office documents, striping is needed. Finally, we intend to develop our scheme for writedominated workloads. The READ strategy in its current form only works well for read-dominated workloads.

Acknowledgments

We thank the anonymous reviewers whose comments noticeably improved the quality of this paper. This work was supported by the US National Science Foundation under grants CNS-0845105 (CAREER), CNS-0834466, and CCF-0742187.

References

- D. Anderson, J. Dykes, E. Riedel, More than an interface—SCSI vs. ATA, in: Proc. USENIX Conf. File and Storage Technologies, 2003, pp. 245–257.
- [2] M. Arlitt, T. Jin, Workload characterization of the 1998 World Cup web site, HP Labs Technical Reports, HPL-1999-35R1, 1999.
- [3] L. Cai, Y.H. Lu, Power reduction of multiple disks using dynamic cache resizing and speed control, in: Proc. Int'l Symp. Low Power Electronics and Design, 2006, pp. 186–190.
- [4] E.V. Carrera, E. Pinheiro, R. Bianchini, Conserving disk energy in network servers, in: Proc. Int'l Conf. Supercomputing, 2003, pp. 86–97.
- [5] D. Colarelli, D. Grunwald, Massive arrays of idle disks for storage achieve, in: Proc. of Supercomputing, 2002, pp. 1-11.
- [6] G. Cole, Estimating drive reliability in desktop computers and consumer electronics systems, Technology Paper TP-338.1, Seagate Technology, November 2000.
- [7] C. Cunha, A. Bestavros, M. Crovella, Characteristics of WWW client-based traces, Technical Report, 1995-010, Boston University, 1995.
- [8] M.D. Dahlin, R.Y. Wang, T.E. Anderson, D.A. Patterson, Co-operative caching: using remote client memory to improve file system performance, in: Proc. USENIX Operating Systems Design and Implementation, vol. 1, 1994. Article No. 19.
- [9] N. Durrant, R. Blish, Semiconductor device reliability failure models, 2000. ismi.sematech.org/docubase/document/3955axfr.pdf.
- [10] J.G. Elerath, Specifying reliability in the disk drive industry: no more MTBF's, in: Proc. IEEE Annual Reliability and Maintainability Symp., January 2000, pp. 194–199.
- [11] Engineering Statistics Handbook, the National Institute of Standards and Technology, NIST, 2007. Website: http://www.sbtionline.com/nist/index. html.
- [12] S. Glassman, A caching relay for the World Wide Web, in: Proc. First Conf. World-Wide Web, 1994, pp. 165–173.
- [13] S. Gurumurthi, A. Sivasubramaniam, M. Kandemir, H. Fanke, DRPM: dynamic speed control for power management in server class disks, in: Proc. Int'l Symp. Computer Architecture, June 2003, pp. 169–179.
- [14] S. Gurumurthi, A. Sivasubramaniam, V.K. Natarajan, Disk drive roadmap from the thermal perspective: a case for dynamic thermal management, in: Proc. Int'l Symp. Computer Architecture, 2005, pp. 38–49.
- [15] S. Gurumurthi, J. Zhang, A. Sivasubramaniam, M. Kandemir, H. Fanke, N. Vijaykrishnan, M. Irwin, Interplay of energy and performance for disk arrays running transaction processing workloads, in: Proc. of ISPASS, March 2003, pp. 123–132.
- [16] G. Herbst, IBM's drive temperature indicator processor (drive-TIP) helps ensure high drive reliability, in: IBM Whitepaper, October 1997.
- [17] Hitachi power & acoustic management—quietly cool, White Paper, Hitachi Corp., March 2004.
- [18] IDEMA Standards, Disk drive reliability benchmark test specification, Document Number R3-98.
- [19] IDEMA Standards, Specification of hard disk drive reliability, Document Number R2-98.
- [20] Y. Kim, S. Gurumurthi, A. Sivasubramaniam, Understanding the performance-temperature interactions in disk I/O of server workloads, in: Proc. 12th Int'l Symp. High-Performance Computer Architecture, 2006, pp. 176–186.
- [21] C.M. Kozierok, Hard disk spindle motor. http://www.pcguide.com/ref/hdd/op/ spin.htm.
- [22] L.W. Lee, P. Scheuermann, R. Vingralek, File assignment in parallel I/O systems with minimal variance of service time, IEEE Trans. Comput. 49 (2) (2000) 127–140.
- [23] Mobile emergency datacenter, North American Access Technologies, Inc., 2006. http://www.naat.com/Disaster%20Recovery/mobile_datacenter.htm.
- [24] A.E. Papathanasiou, M.L. Scott, Power-efficient server-class performance from arrays of laptop disks, in: WIP Presentation at the USENIX Annual Technical Conf., Boston, June 27–July 2, 2004.
- [25] D.A. Patterson, G. Gibson, R.H. Katz, A case for redundant arrays of inexpensive disks (RAID), in: Proc. ACM SIGMOD Int'l Conf. Management of Data, 1988, pp. 109–116.
- [26] E. Pinheiro, R. Bianchini, Energy conservation techniques for disk array-based servers, in: Proc. ACM Int'l Conf. on Supercomputing, June 2004, pp. 68–78.
- [27] E. Pinheiro, W. Weber, L. Barroso, Failure trends in a large disk drive population, in: Proc. 5th USENIX Conf. File and Storage Technologies, San Jose, CA, February, 2007.
- [28] Power, heat and sledgehammer, White Paper, Maximum Institution Inc., April 2002.
- [29] Reliability of hard disk drives (HDD), Technical Paper, University of Maryland, CALCE, 2003.

- [30] B. Schroeder, G.A. Gibson, Disk failures in the real world: what does an MTTF of 1,000,000 hours mean to you? in: Proc. 5th USENIX Conf. File and Storage Technologies, San Jose, CA, 2007.
- [31] T. Schwartz, M. Baker, S. Bassi, B. Baumgart, W. Flagg, C. Ingen, K. Joste, M. Manasse, M. Shah, Disk failure investigations at the internet archive, in: Proc. 14th NASA Goddard, 23rd IEEE Conf. Mass Storage Systems and Technologies, May 2006.
- [32] S. Shah, J.G. Elerath, Reliability analysis of disk drive failure mechanisms, in: Proc. IEEE Reliability and Maintainability Symp., 2005, pp. 226–231.
- [33] S. Sobti, N. Garg, C. Zhang, X. Yu, A. Krishnamurthy, R.Y. Wang, PersonalRAID: mobile storage for distributed and disconnected computers, in: Proc. First Conference on File and Storage Technologies, FAST, January 2002.
- [34] S.W. Son, G. Chen, M. Kandemir, Disk layout optimization for reducing energy consumption, in: Proc. 19th Annual Int'l Conf. Supercomputing, 2005, pp. 274–283.
- [35] S.W. Son, M. Kandemir, Energy-aware data prefetching for multi-speed disks, in: Proc. ACM Int'l Conf. Computing Frontiers, Ischia, Italy, May 2006.
- [36] S.W. Son, M. Kandemir, A. Choudhary, Software-directed disk power management for scientific applications, in: Proc. Int'l Symp. Parallel and Distributed Processing, April 2005.
- [37] C. Weddle, M. Oldham, J. Qian, A. Wang, PARAID: the gear-shifting poweraware RAID, Technical Report 060323, Florida State University, January 2006.
- [38] Q. Xin, E.L. Miller, T. Schwarz, D.D.E. Long, S.A. Brandt, W. Litwin, Reliability mechanisms for very large storage systems, in: Proc. 20th IEEE Conf. Mass Storage Systems and Technologies, 2003, pp. 146–156.
- [39] H. Yada, H. Ishioka, T. Yamakoshi, Y. Onuki, Y. Shimano, M. Uchida, H. Kanno, N. Hayashi, Head positioning servo and data channel for HDDs with multiple spindle speeds, IEEE Trans. Magn. 36 (5) (2000) 2213–2215.
- [40] J. Yang, F. Sun, A comprehensive review of hard-disk drive reliability, in: Proc. IEEE Reliability and Maintainability Symp., January 1999, pp. 403–409.
- [41] S. Yin, X. Ruan, A. Manzanares, X. Qin, How reliable are parallel disk systems when energy-saving schemes are involved? in: Proc. IEEE International Conference on Cluster Computing, CLUSTER, New Orleans, LA, August 31–September 4, 2009.

- [42] Q. Zhu, Z. Chen, L. Tan, Y. Zhou, K. Keeton, J. Wilkes, Hibernator: helping disk arrays sleep through the winter, in: Proc. 12th ACM Symp. Operating Systems Principles, 2005, pp. 177–190.
- [43] Q. Zhu, Y. Zhou, Power-aware storage cache management, IEEE Trans. Comput. 54 (5) (2005) 587–602.



Tao Xie received the Ph.D. degree in computer science from the New Mexico Institute of Mining and Technology in 2006. He received the B.Sc. and M.Sc. degrees from Hefei University of Technology, China, in 1991 and 2000, respectively. He is currently an assistant professor in the Department of Computer Science at San Diego State University, San Diego, California. He received a US National Science Foundation (NSF) Early Career Award in 2009. His research interests are storage systems, security-aware scheduling, high performance computing, cluster and Grid computing, parallel and distributed systems, and real-

time/embedded systems. He is a member of the IEEE.



Yao Sun received the master degree in computer science from San Diego State University in 2008 with the honor "the most outstanding graduate student". He received his B.S. degree in computer science and engineering from Harbin Institution of Technology, China in 2002. He is currently a software engineer at Teradata, which offers the industry's only real distribute computing database product.