

Performance Analysis of Disk Arrays Under Failure*

Richard R. Muntz and John C.S. Lui

UCLA Computer Science Department, LA, CA 90024-1596, USA

Abstract

Disk arrays (RAID) have been proposed as a possible approach to solving the emerging I/O bottleneck problem. The performance of a RAID system when all disks are operational and the $MTTF_{sys}$ (mean time to system failure) have been well studied. However, the performance of disk arrays in the presence of failed disks has not received much attention. The same techniques that provide the storage efficient redundancy of a RAID system can also result in a significant performance hit when a single disk fails. This is of importance since single disk failures are expected to be relatively frequent in a system with a large number of disks. In this paper we propose a new variation of the RAID organization that has significant advantages in both reducing the magnitude of the performance degradation when there is a single failure and can also reduce the $MTTF_{sys}$. We also discuss several strategies that can be implemented to speed the rebuild of the failed disk and thus increase the $MTTF_{sys}$. The efficacy of these strategies is shown to require the improved properties of the new RAID organization. An analysis is carried out to quantify the tradeoffs.

1 Introduction.

A disk array is a set of disks with redundancy to protect against data loss. Patterson [7] describes sev-

*This research was supported by IBM through a University of California MICRO grant.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.

eral methods of organizing the data on disks including mirroring and multiple data blocks plus parity. Compared with mirroring, the " $N + 1$ " RAID organization sacrifices some performance in terms of I/O access rate but with a reduced storage penalty for redundancy. In the remainder of this paper we will use the term RAID to refer to the $N + 1$ RAID organization. We also concentrate exclusively on the "small read/small write" use of the disks in which they are accessed independently rather than the mode in which large reads (writes) are performed by concurrent reads (writes) executed on each disk.

Two aspects of RAID systems have been examined most closely: performance under the condition that all disks are operational (e.g., [1]) and the mean time to system failure (data loss) (e.g. [4]). The performance of RAID when there are inoperable (or inaccessible) disks has largely been ignored (an exception is [8]).

Compared to system failure, it is a relatively common occurrence to have a single disk unavailable and therefore the performance of the system during the repair period is of concern. The $N + 1$ RAID organization achieves a low cost in redundant storage overhead but at the price of requiring multiple reads to the surviving disks in the same array each time a block on the failed disk must be reconstructed (i.e., to satisfy a read request for that block). In the worst case (a workload of all reads and no writes) this can double the access rate to the surviving disks and thus in effect, cut the capacity of the array in half. Consider for example a shared nothing DBM architecture as in Figure 1. Each node in this system would have one or more disk arrays. The impact that this can have on total system performance is dependent on the characteristics of the system workload. We are particularly interested in database applications of the RAID architecture. The impact of a single disk failure is most severe in the case of a "decision support" environment in which complex queries are common and the database tables have been partitioned among the disks on all or many nodes for increased I/O bandwidth. Complex operations can be limited by any imbalance in the system, which can be caused

by skew in the load [6] or by a disk array with diminished capacity due to a failed disk [5]. In a one hundred disk system, a single failed disk will represent a loss of only 1% of the raw I/O capacity of the system. However, if the effect is to reduce the capacity of the array to which it belongs by say 25%, this can cause a significant imbalance in the system, and the impact on aggregate system performance can be considerable.

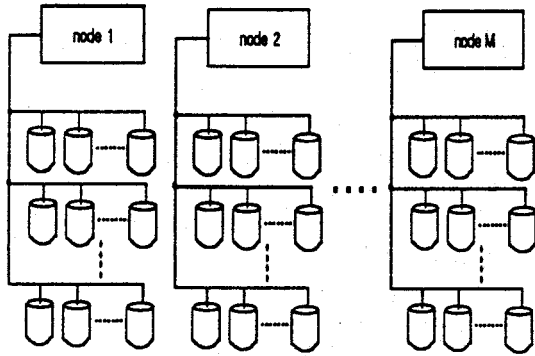


Figure 1: Shared Nothing DBM with disk arrays.

In order to maintain a reasonable $MTTF_{sys}$ (mean time to failure of the system) in a system with a large number of disks it is necessary to provide immediate repair of a failed disk. To accomplish this, "hot standby" disks are provided and the system automatically rebuilds the contents of the failed disk on the standby disk from the redundant information on the surviving disks [4]. The $MTTF_{sys}$ of the RAID is easily shown to be inversely proportional to the rebuild time [1, 3]. In the $N+1$ RAID system described in [7], to rebuild the failed disk contents at maximum speed (the capacity of the standby disk) would use the entire capacity of the surviving disks in the array. Thus to rebuild at maximum rate would mean that the array can perform no other work during the rebuild period. One can of course, tradeoff the rebuild rate with the rate at which the surviving disks process normal workload requests. However, this increases the time to rebuild the failed disk contents and thereby decrease the $MTTF_{sys}$. Quantifying this tradeoff is one of the purposes of this paper.

We also propose a new RAID organization. In the proposed architecture the same type of tradeoffs exist between normal workload and $MTTF_{sys}$, but the proposed organization can (for some system and workload parameters) yield an improved $MTTF_{sys}$ and support a higher workload during the rebuild period.

This is accomplished by relaxing an assumption that has been made in most of the RAID work, i.e., that the "group size" (the number of data blocks plus parity) is the same as the "cluster" size (the number of disks over which the groups of blocks are distributed). We propose consideration of larger cluster sizes and show that by spreading the groups over a larger cluster of disks the increased load (per disk) on surviving disks can be significantly decreased. The cluster size will effect the time required to rebuild a failed disk (and therefore the $MTTF_{sys}$) and also the workload (measured in accesses per second per disk) that can be supported during the rebuild. These issues are addressed analytically in a later section. The results indicate that substantially better performance and reliability can often be obtained from such an organization with properly chosen values for group and cluster size.

The only study that we are aware of that considers disk subsystem performance under failure and the reconstruction of the failed disk is the Copeland and Keller study comparing mirroring and "interleaved declustering" [3]. The analysis of the proposed RAID organization presented here is similar with the methodology introduced in [3].

In section 2 we present the proposed new RAID organization and also several strategies for efficient rebuilding of a failed disk. In section 3 we present an outline of the analysis of the proposed techniques. Section 4 presents numerical results for a range of parameter values and their interpretation. Section 5 summarizes the contributions of the paper.

2 RAID Recovery.

It is assumed here (as in [3]) that there is a system requirement to provide some minimum level of service during the rebuild period which will be measured in terms of accesses per second per disk, denoted by λ_0 . λ_0 is the access rate per "logical disk" and when an array contains a failed disk, the load on the surviving disks will actually be higher since the λ_0 accesses per second originally directed to the failed disk will require additional I/Os to the surviving disks. A consequence of the requirement for a minimum workload throughput is an additional limitation on the rate at which the failed disk can be rebuilt, i.e. the full capacity of the surviving disks is not available for rebuilding the contents of the failed disk. The cost of reconstructing the contents of a block from the failed disk causes the tradeoff between λ_0 and the rebuild rate to be of greater importance because the rebuild time becomes more sensitive to λ_0 . Next we introduce

for a data block) the addresses of the buddy blocks must be either easily computable or the set of different groups should be small enough to allow for table lookup. This is still an open problem which we are currently investigating. It appears to be a combinatorial block design problem [2]. This paper addresses the issue of evaluating the proposed architecture under the assumption that the group assignment and addressing issues will be adequately solved given that we can show sufficient advantages to the larger cluster size.

2.2 Disk Rebuild Strategies.

In this section we describe several options that might be employed in rebuilding the failed disk. We start with a description of the simple baseline copy procedure and then discuss the more sophisticated schemes in following subsections.

An assumption common in all the rebuild schemes we consider concerns how writes to the failed disk are handled. It is clear that writes of blocks that have already been copied to the standby disk should be redirected to the standby to keep that portion of the disk up to date. For writes to the portion of the disk that has not yet been copied there are two possible options: (1) convert the write to a null operation or, (2) write the block to the corresponding block on the standby disk. We will assume in all schemes that all writes to the failed disk are redirected to the standby disk. This appears to be relatively easy to accomplish using a bit map to indicate which blocks have already been copied. The copy procedure will just skip over the blocks that have already been copied as a result of the write operations. The bit map would require no more than 16 Kbytes for reasonable block size and disk capacity.

Baseline Copy Procedure.

The baseline procedure then simply sequentially reads blocks from the failed disk (causing reconstruction) and writes them to the standby disk. The only exception to the sequential copying is to skip over blocks already copied due to writes in the normal workload.

Rebuild with Redirection of Reads.

When the failed disk has been partially reconstructed on the standby, it is possible to satisfy some reads by either reconstructing the block (by multiple reads of blocks on the surviving disks) or by reading from the standby (if the block has already been copied). Which of the two options is optimal will depend on whether the surviving disks are the bottleneck in the rebuild process or the standby disk is

the bottleneck (and the optimum choice can change with time). Let $f(t)$ be the fraction of the failed disk that has been reconstructed at time t (the rebuild starts at $t = 0$). If we assume that the access pattern of the normal workload is uniform across all blocks, then the fraction of reads that can be redirected to the standby disk at time t can be denoted by $b(t) \cdot f(t)$ where $b(t)$ is a control variable which can vary between 0 and 1. In the analysis presented in the next section we determine the optimal policy and assume perfect adherence to that optimal policy.

Piggy-backing Rebuild on Normal Workload.

The idea here is to "capture" a block from the failed disk that is reconstructed due to a read request that was issued as part of the normal workload. A relatively simple implementation is possible assuming only a slight modification to a conventional buffer manager. When a read from the failed disk is satisfied by reconstruction of the block, the block contents will be stored in a buffer. At this time buffer copy can be marked as "modified", the disk address associated with the buffer page can be changed to refer to the standby disk, and the rebuild bit map can be changed to show this page as being copied. Normal buffer replacement will eventually copy the contents to the appropriate block on the standby. (At the end of the rebuild any blocks remaining in the cache can be flushed to the standby disk.)

Piggy-backing and redirection can help to reduce the load on the surviving disks in an array which has a failed disk. The extent to which this can be of benefit depends on several factors. For the same rebuild rate and the same minimum access rate for the normal workload, both redirection and piggy-backing can be used to reduce the load on the surviving disks. The I/O access capacity thus made available on the surviving disks can be used to either support a higher normal workload during rebuild (λ_0) or a faster rebuild. In the following section we present an analysis of disk arrays under failure to quantify the tradeoffs between (1) the cluster size C , (2) the minimum workload that must be supported λ_0 and, (3) the $MTTF_{sys}$. Qualitatively the tradeoffs are clear but quantitatively they are not. With a constant minimum workload λ_0 , a larger cluster size will allow a faster copy rate but the larger cluster size increases the likelihood of a second failure in the cluster. Which effect is dominant is not immediately obvious. As we will show from the analysis in the next section, a larger cluster size can sometimes simultaneously support a higher minimum workload and provide a longer $MMTF_{sys}$.

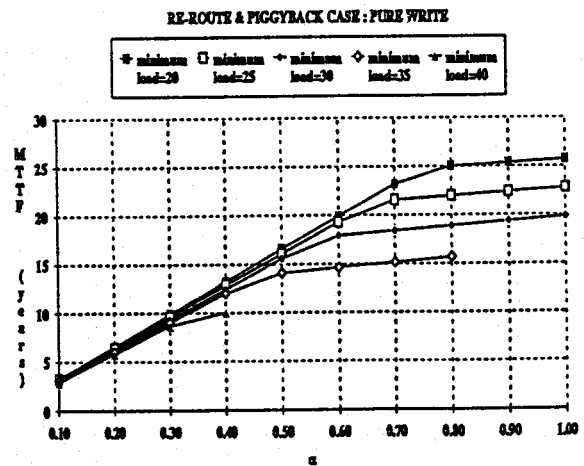
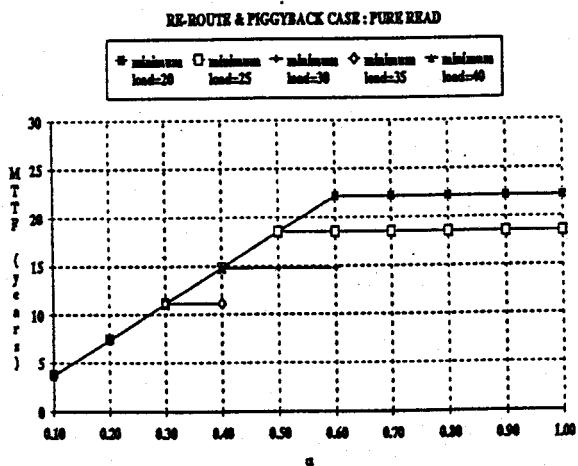
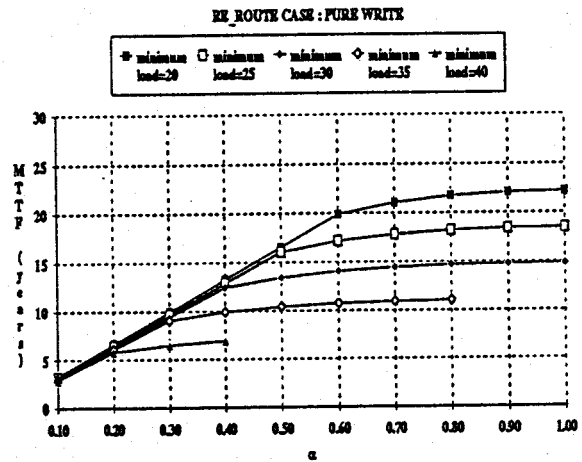
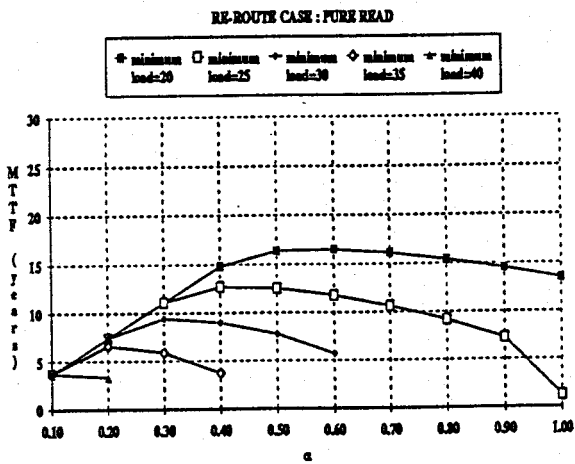
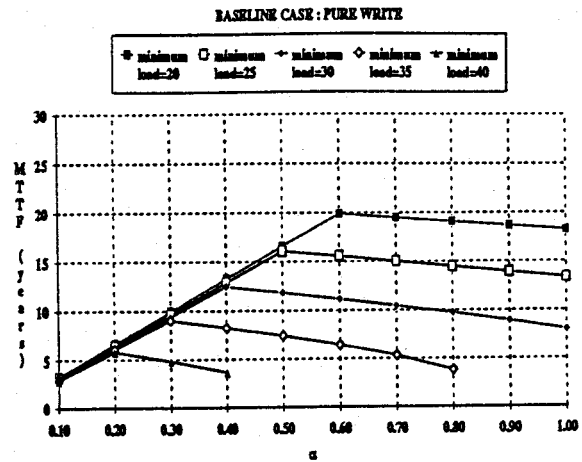
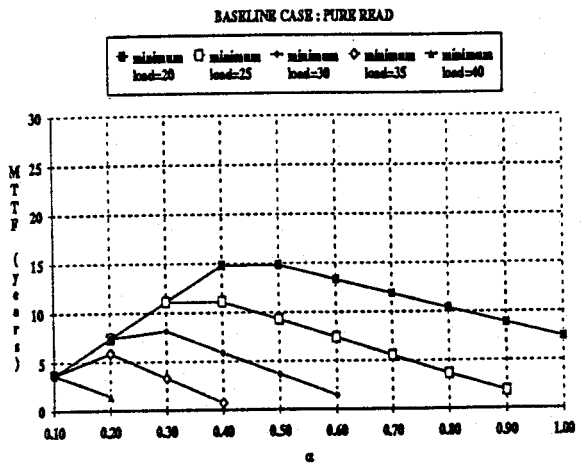


Figure 7: $MTTF_{sys}$ for various recovery scheme for pure read transactions.

Figure 8: $MTTF_{sys}$ for various recovery scheme pure write transactions.

The rebuild time T will be :

$$T = t_1 + \frac{I}{\lambda_0(\frac{\rho}{2} - 1)} \ln \left[\frac{C_1 - 1 + f(t_1)}{C_2} \right] \quad (16)$$

where :

$$C_1 = \frac{2K}{\lambda_0(2-\rho)} + \frac{\rho[1-f(t_1)]}{2-\rho} + 1 - f(t_1)$$

$$C_2 = \frac{2K}{\lambda_0(2-\rho)} + \frac{2[1-f(t_1)]}{2-\rho}$$

Case 3: if $\frac{\mu(1-\alpha)}{1+\alpha(\rho-1)} \leq \lambda_0 \leq \frac{\mu}{1+\alpha\rho}$.
for $t^* \leq t_1$:

$$\lambda_c^*(t) = \frac{\mu - \lambda_0}{\alpha(1-\rho)} - \left[\frac{(\mu - \lambda_0)\rho}{\alpha(1-\rho)} + \rho\lambda_0 \right] e^{-\frac{\lambda_0}{2}(1-\rho)t} \quad (17)$$

for $t^* > t_1$:

$$\lambda_c^*(t) = \left(\frac{2K + \lambda_0\rho[1-f(t_1)]}{2-\rho} \right) - \left(\frac{K\rho + \lambda_0\rho[1-f(t_1)]}{2-\rho} \right) e^{-\frac{\lambda_0}{2}(1-\rho/2)(t-t_1)} \quad (18)$$

where :

$$f(t_1) = \frac{\mu - \lambda_0 + \alpha\lambda_0(1-\rho)}{\alpha\lambda_0(1-\rho)} \left[1 - e^{-\frac{\lambda_0}{2}(1-\rho)t_1} \right]$$

$$K = (\mu - \lambda_0)\left(1 - \frac{\alpha-1}{2\alpha}\right) - \frac{1}{2}\lambda_0\rho + \frac{1}{2}\rho\lambda_0 f(t_1)$$

$$t^* = -\frac{I}{\lambda_0(1-\rho)} \ln \left[1 - \frac{\alpha\lambda_0(1-\rho)}{\mu - \lambda_0 + \alpha\lambda_0(1-\rho)} \right]$$

$$t_1 = -\frac{I}{\lambda_0 w} \ln \left[\frac{(\mu - \lambda_0)(1-\alpha w)}{(\mu - \lambda_0)\rho + \alpha\rho\lambda_0 w} \right]$$

The rebuild time T will be :

If $t^* \leq t_1$:

$$T = t^* \quad (19)$$

If $t^* > t_1$:

$$T = t_1 + \frac{I}{\lambda_0(\rho/2 - 1)} \ln \left[\frac{C_1 - 1 + f(t_1)}{C_2} \right] \quad (20)$$

where :

$$C_1 = \frac{2K}{\lambda_0(2-\rho)} + \frac{\rho[1-f(t_1)]}{2-\rho} + 1 - f(t_1)$$

$$C_2 = \frac{2K}{\lambda_0(2-\rho)} + \frac{2[1-f(t_1)]}{2-\rho}$$