

## Two Papers on RAIDs

*Peter Chen, Garth Gibson, Randy H. Katz  
David A. Patterson, Martin Schulze*

Computer Science Division  
Department of Electrical Engineering and Computer Sciences  
University of California  
Berkeley, CA 94720

### *ABSTRACT*

A RAID is a Redundant Array of Inexpensive Disk, a new way to organize small format disk devices to drastically increase I/O bandwidth.

In this technical report, we describe the RAID concept, the basic RAID levels, a more detailed analysis of RAID performance on reliability, and our initial prototyping plans.



# Introduction to Redundant Arrays of Inexpensive Disks (RAID)

*David A. Patterson, Peter Chen, Garth Gibson, and Randy H. Katz*

Computer Science Division  
Department of Electrical Engineering and Computer Sciences  
571 Evans Hall  
University of California  
Berkeley, CA 94720  
(patt@cs.berkeley.edu)

## 1. The Pending I/O Crisis

The computer industry has entered a period of unprecedented improvement in CPU performance. Midrange uniprocessors are improving at a rate of 50% to 100% per year, with this uniprocessor rate multiplied by the increasing popularity of multiprocessors.

There is more to a computer system, however, than the processor. Memory and I/O must match these gains to deliver a system that achieves the potential of the processor. Various rules of thumb tie CPU performance to both the capacity and the speed of the memory and I/O subsystems, so we need advances in both dimensions to create faster yet balanced computers.

Memory subsystems are matching the challenge of CPU performance. DRAMs are growing in capacity by about the same rate as CPUs are improving in performance. The performance of DRAMs is improving at a much more modest rate, perhaps doubling every decade. SRAMs, however, are matching the performance improvement of CPUs. Fortunately there is a long list of architectural innovations--duplicated caches, cache coherency, multilevel caches, prefetching, interleaved memory, pipelined memory, and so on--allowing fast SRAMs and large DRAMs to be combined into memory systems that can match the performance demands of new processors.

I/O systems performance is limited by networks and magnetic disks. There are several efforts to improve network speeds by factors of 10 to 100, so we see little trouble here provided some of these efforts succeed. The good news about magnetic disks is that improvements in capacity and cost per megabyte are keeping pace with processors. The bad news is that performance gains are modest. Rotation speed is unchanged in more than a decade, and in that same time period seek time has improved by no more than a factor of two.

Without innovation, we see most programs becoming I/O bound. If such an I/O crisis comes to pass, there will be little reason to buy faster processors, since it is economic nonsense to pay more to increase processor idle time.

## 2. Arrays of Inexpensive Disks

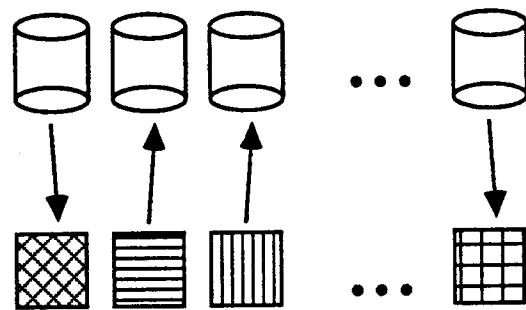
While the magnetic disk industry has made little progress in improving speed of disks, it has significantly reduced the size of disks. The personal computer industry has created a market for 5.25 and 3.5 inch drives, reducing the cost per disk system as well as the traditional lowering of cost per megabyte. Table I below compares the top-of-the-line IBM 3380 model AK4 mainframe disk, Fujitsu M2361A "Super Eagle" minicomputer disk, Impres/CDC Wren-IV workstation disk, and the Conner Peripherals CP 3100 personal computer disk.

<i>Characteristics</i>	<i>IBM 3380</i>	<i>Fujitsu M2361A</i>	<i>CDC Wren-IV</i>	<i>Conners CP3100</i>
Disk diameter (inches)	14	10.5	5.25	3.5
Formatted Data Capacity (MB)	7,500	600	340	100
MTBF rated by manufacturer (hours)	30,000	20,000	40,000	30,000
No. Actuators	4	1	1	1
Maximum I/O's/second/Actuator	50	40	40	30
Maximum I/O's/second/box	200	40	40	30
Transfer Rate (MB/sec)	3	2.5	1.5	1
Power/box (W)	1,650	640	40	10
MB/W	4.5	0.9	8.5	10.0
Volume (cu. ft.)	24	3.4	0.3	.1
MB/cu. ft.	312	176	1133	1000

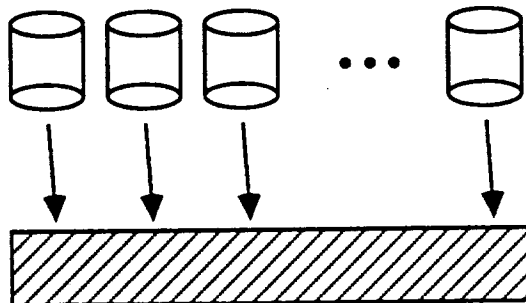
**Table I.** Comparison of IBM 3380 disk model AK4 for mainframe computers, the Fujitsu M2361A "Super Eagle" disk for minicomputers, Impress/CDC Wren-IV disk for workstations, and the Conners Peripherals CP 3100 disk for personal computers. By "Maximum I/O's/second" we mean the maximum number of average seeks and average rotates for a single sector access.

The table shows that the small drives are close to the performance and the reliability of the large drives. The biggest difference is the lower power and smaller volume per megabyte, with the smallest drive about three to five times better. What is not listed is cost per megabyte. The trade-off is the potential lower cost per megabyte of large drives--since they can amortize the cost of the support electronics over a larger number of megabytes--versus the economies of scale provided by the higher sales volume of the smaller drives. It is also an issue of manufacturer's cost versus selling price, with the likelihood that the mainframe drives having larger markup. In our discussions with companies that make several sizes of disks they say the cost per megabyte is a function of the popularity of a given disk at a given time; they suggest assuming the cost per megabyte is independent of disk diameter. We make that assumption.

Given similar performance and cost per megabyte, one way to address higher performance is replace a single large drive by an array of many smaller drives.[Kurzweil 88] Such an array provides many more arms per dollar, meaning higher performance because many small request can be serviced independently *and* large requests can be spread over several disks to transfer in parallel (Figure 1). Moreover, the advantages in volume and power can mean a smaller footprint and lower air conditioning requirements.



(a) Arrays of inexpensive disks support several small, independent reads and writes simultaneously.



(b) Arrays of inexpensive disks also support large reads or writes, where all disks transfer in parallel.

**Figure 1.** Replacing a single large expensive disk by an array of inexpensive disks improves performance because it can support (a) many small individual accesses simultaneously and (b) large accesses with all disks transferring in parallel ("striping" [Kim 86] [Salem 86]).

### 3. Redundant Arrays of Inexpensive Disks

The drawback to replacing a single large disk with, say, 100 small disks is reliability. Basically 100 devices have 1/100th the reliability of a single device, reducing the mean time between failure (MTBF) from over three years to less than two weeks. This is so poor that without a scheme to improve reliability, arrays containing many disks are unfeasible.

Fortunately redundancy can improve reliability<sup>1</sup> of 100 small disks beyond that of a single large disk. Although failures occur 100 times more frequently with 100 disks, the chances of a second failure before the first is replaced is small enough to tolerate more failures and still be more reliable than a single disk. Thus Redundant Arrays of Inexpensive Disks, or RAID, has the potential advantage of not only higher performance with lower power and smaller footprint, but also higher reliability.

In an earlier paper we presented five different schemes for disk redundancy [Patterson 88], but here we only present the two schemes most likely to be implemented. If we include the rest of the computer system it would seem we would need to duplicate all components to achieve high reliability. A companion paper shows high reliability is possible to achieve with more modest redundancy costs [Schulze 89]. In this paper we assume that the reliability is sufficient. Readers interested in redundancy schemes for much larger disk arrays, file systems or databases for RAID should see [Gibson 89], [Douglass 89], or [Stonebraker 88].

1. In this paper we use the term reliability to include availability.

## 4. Mirrored RAID

The simplest redundancy scheme is to double the number of disks, keeping a redundant copy of each datum. If a disk fails, the system uses the redundant copy until the failed disk is replaced, and then copies a redundant version to the new disk. Data is lost only if the other disk of the pair fails before first is replaced. In normal operation a copy is maintained by making every write update both disks. This scheme is variously called *mirroring*, *shadowing*, or *copying*, and is used by Tandem, DEC, and IBM to improve reliability. (In our original paper we called mirroring a level 1 RAID.)

This scheme has the highest cost: the user must double number of disks for the same amount of data or, conversely, use only half the real storage capacity of the disks. If the arms and spindles of a pair were synchronized then the performance of mirroring versus nonredundant disks would be the same. This is not commonly how the mirroring is implemented, and a write results in independent writes to two disks. They can be overlapped, but in general one will have longer seek and/or rotational delay. On the other hand, the independence of disks can improve performance of reads. The system might look at the pair of disks that have the data; if only one is busy, it chooses the other. If both are idle, it picks the disk that has shortest seek [Bitton 88].

In summary, mirrored RAID's have the highest cost for a given storage capacity, but performance versus a nonredundant disk array depends on the mix of reads and writes.

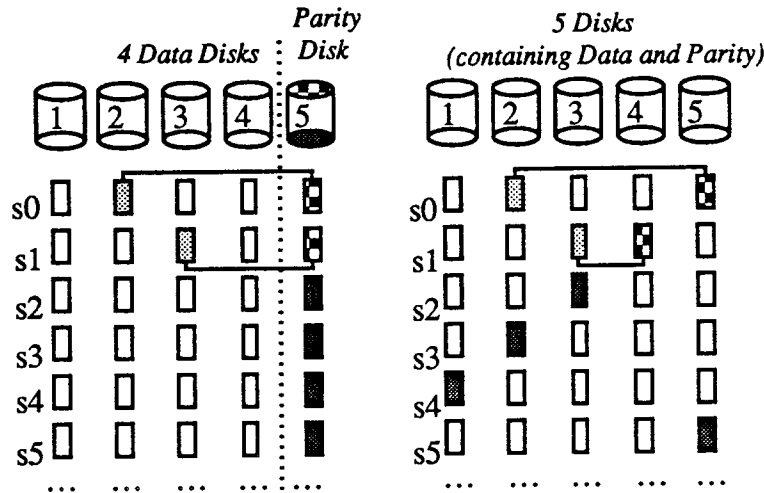
## 5. N+1 RAID

An advantage of disks is that they can detect their own mistakes: either the disk controller will not get a response or the ECC code per sector will be incorrect. By calculating and storing parity of a group of disks on a bit-per-disk basis, any single disk failure can be corrected simply by reading the rest of the disks in the group to determine what bit value on the failed disk would give the proper parity. This N+1 RAID can lose data only if there is a second failure in the group before the failed drive is replaced.

This scheme has much lower cost overhead, with the customer deciding how much overhead he wants to pay by increasing the number of disks in the parity group. Performance depends not only on the mix of reads and writes, but also on the size of the accesses. Since there is ECC information on each sector, read performance is essentially the same as nonredundant disk arrays. For "large" writes--writing at least a sector to every disk in the parity group--the only performance hit is 1/N more writes to write the parity information. Writes to data on a single disk, on the other hand, require four disk accesses:

- 1) Read the old data;
- 2) Read the old parity;
- 3) Write the new data;
- 4) Write the new parity using this formula:  
$$\text{new parity} = (\text{old data} \text{ xor } \text{new data}) \text{ xor } \text{old parity}$$

It would seem that an additional performance limit would be the parity disk, since small writes to any disk must also cause a read and a write to the parity information. Such a bottleneck is avoided by spreading the parity over several disks. Figure 2 shows how the straight-forward implementation is altered to avoid parity bottlenecks. (In our original paper we called N+1 a level 5 RAID.)



(a) Parity information in straight-forward scheme. The sectors are shown below the disks. (The checked areas indicate the parity information.) Writes to s0 of disk 2 and s1 of disk 3 imply writes to s0 and s1 of disk 5. The parity disk (5) becomes the write bottleneck.

(b) Parity information for rotated scheme. The sectors are shown below the disks, with the parity information and data spread evenly through all the disks. Writes to s0 of disk 2 and s1 of disk 3 still imply 2 writes, but they can be split across 2 disks: to s0 of disk 2 and to s1 of disk 4.

**Figure 2.** The performance impact of this small change is large in large parity groups since it allows N+1 RAID to support multiple individual writes per group. For example, suppose we want to write sector 0 of disk 2 and sector 1 of disk 3. As shown on the left writes must be sequential since both sector 0 and sector 1 of disk 5 must be written. However, as shown on the right, writes can proceed in parallel since a write to sector 0 of disk 2 still involves a write to disk 5 but a write to sector 1 of disk 3 involves a write to disk 4.

## 6. Performance of Mirrored RAID vs. N+1 RAID

Comparing these two RAID organizations is both simple and difficult. Common sense suggests Mirrored RAID, using roughly twice as many disks, is more expensive and has higher performance. If cost is your only concern, you pick N+1, and you pick Mirroring if performance is the only concern. What if you care about both cost and performance? We use the metric of throughput per disk, since a customer can always buy more disks to solve an I/O bottleneck in any scheme. From the discussion above it is clear that comparative performance is then sensitive to whether the accesses are reads or writes and the size of the accesses.

By making several simplifying assumptions we can estimate comparative performance in all these measures. These simplifying assumptions include:

- Every access is assumed to take one average seek and one average rotation;
- Every access is the same size;
- Accesses are spread optimally across all disks;
- Disks are never idle waiting for requests;
- Accesses are assumed to be homogeneous, e.g., 100% small reads.
- There is no optimization to schedule reads on mirrored disks;
- Latency is ignored.

Figure 3 is a comparison of the two schemes using these assumptions, with 100% being the performance of a nonredundant disk array for that type of access. Using these assumptions we see that read performance is identical, with N+1 winning on capacity and large writes while Mirroring wins on small writes.



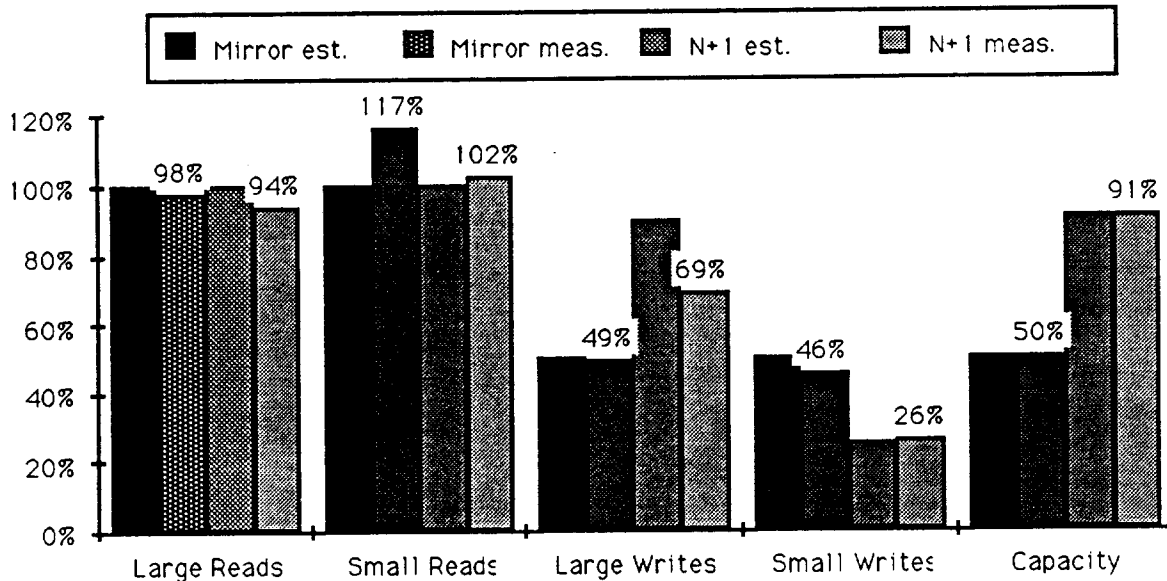
**Figure 3.** Estimated performance of Mirrored RAID vs. N+1 RAID for Large Reads, Small Reads, Large Writes, Small Writes, and Useful Storage Capacity. Small accesses mean one block while large accesses mean one block per disk in the parity group. The group size for N+1 is 11 disks. Measures are percentage of nonredundant disk throughput.

To see if these results would hold in more realistic circumstances, we recently completed an experiment on an Amdahl 5890 using many large Amdahl 6380 devices[Chen 89]. N+1 RAID used 11 6380 devices with a track of 4 KB as the primary block parity block size. This experiment improved the above calculation in the following ways:

- Real hardware was used, accounting for CPU time and xor calculation time;
- Seek and rotation times were not constant;
- The size of large accesses is not exactly one block per disk in a parity group and a large access is not aligned to fit in the minimum number of parity groups (see Figure 5 below);
- Accesses varied in size around an average. The average size was 6.5 KB for small accesses and 1.5 MB for large accesses. Several distributions of access sizes were used;
- Accesses were not spread evenly between all disks, so some were hot spots and some were underutilized;
- Latency is considered, with the the load varied until 90% of accesses met a latency threshold;
- For mirrored disks the reads were optimized to minimize seeks, thereby slowing mirrored writes.

Figure 4 compares the measured results versus the estimated results in Figure 3.



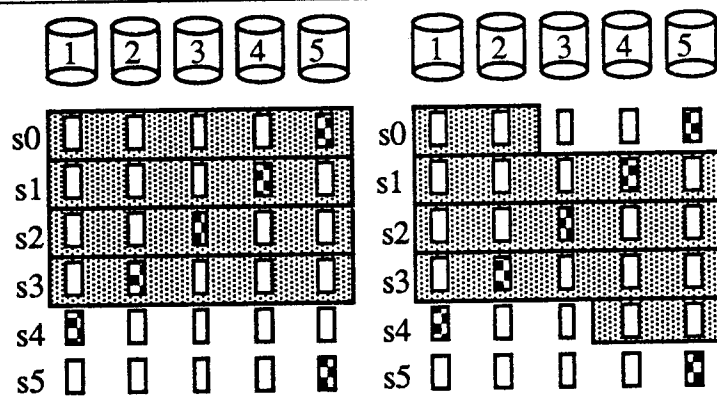


**Figure 4.** Estimated and measured performance of Mirrored RAID vs. N+1 RAID for Large Reads, Small Reads, Large Writes, Small Writes, and Useful Storage Capacity. The experiment was run on an Amdahl 5890 CPU using Amdahl 6380 devices. For N+1 the group size was 11. The average size of a small access is 6 KB and the average size of a large access is 1.5 MB. Measures are percentage of nonredundant disk throughput.

While most of the measured results were close to the estimates, a few were not:

- Large reads for N+1 RAID did not achieve 100% because of the parity information. While parity need not be read, it still takes time for sector containing the parity to spin under the head or for the head to move over the track containing parity during a seek. For the large writes in this experiment about 3 of the 44 tracks would contain parity, and 41/44 ratio of data to total space is close to 94% of total time that we measured.
- The optimization choosing which mirrored disk was best to minimize seek distance of small reads over a nonredundant disk by 17%. This optimization had little effect on large transfers.
- N+1 RAID large writes were 69% vs. 91% of nonredundant disks because the accesses were not "aligned" to exactly one track per disk. Figure 5 shows the model for the estimate and the more realistic model for the measurements. An average large write in this experiment would write 75% of the tracks in a full parity group write with the remaining 25% split across two partial parity group writes.
- Small writes were slightly faster for N+1 RAID and somewhat slower than expected for Mirrored RAID. Small writes do not need to seek to write the new values after reading the old in N+1, they just pay a full rotation waiting for the old data to spin under the head again. Mirrored RAID small writes were slower than expected, in part because of the waiting for the longer of the two rotational delays.

Using the assumptions and the results from this experiment we see a closing of the performance gap on writes--Mirroring is closer to N+1 on large writes and N+1 is slightly closer to Mirroring on small writes--while Mirroring gains a slight edge on reads, with capacity still on the side on N+1.



(a) In estimate we assumed that a large writes were multiples of the size of the parity group, and that they were aligned so that there were no partial groups.

(b) The writes to s0 of disks 1 and 2 will require reads in the s0 parity group to calculate the parity in s0 of disk 5, and similarly extra reads will be needed for parity group s4.

**Figure 5.** The estimates used the model on the left (a), while the experiment randomly placed the large files so there would be writes to portions of two parity groups in addition to the writes to full parity groups. Clearly the importance of alignment is dependent on the size of a large write.

## 7. RAID-I Prototype

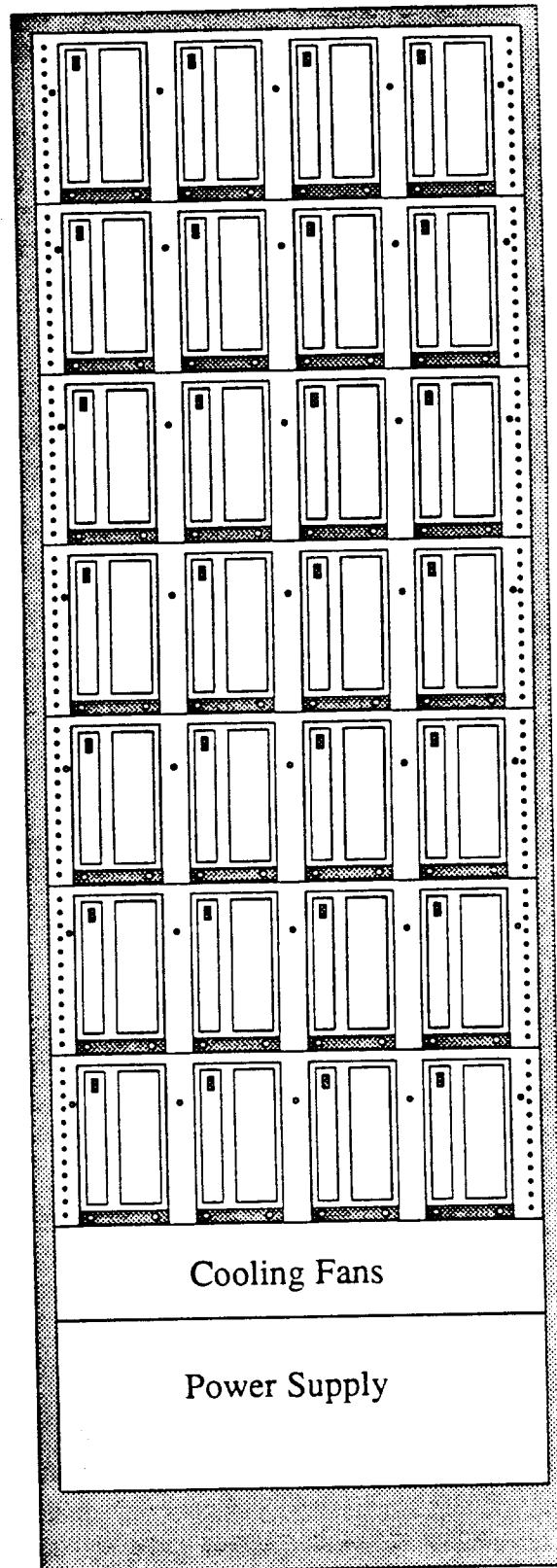
The RAID hardware research is not being done in isolation. The XPRS project--eXperimental Postgres, Raid, and Sprite--includes the development of the Sprite operating system and the Postgres database to take advantage of higher performance I/O systems. XPRS is in turn part of Mammoth project of U.C. Berkeley Computer Science Division that is exploring the advantages of massive storage across many fields of computer science.

To provide a vehicle for architectural experiments and software development for the XPRS and Mammoth projects, we are constructing a prototype we call RAID-I. It consists of:

- Sun-4/280 with 128 MB of memory;
- 7 32-bit VME-SCSI Host/Bus Adapters;
- 32 CDC Wren-IV 340 MB 5.25" disks (using embedded SCSI controllers);
- 1 Ethernet interface.

Figure 6 shows drawing of RAID-I.

RAID-I is an off-the-shelf system, to be used for software development and for measurements to determine the design of later RAID systems. Hence RAID-I has no performance goals, with the only goal being stability for software development. We expect that the VME bus, operating system overhead, and SCSI controller overhead will conspire to reduce the available bandwidth to much below that of 32 disks. Our goal is for RAID-I to have enough capacity to make it an attractive resource in our department so that we can get hands on experience with RAIDs with real users. We expect RAID-I to be running Sprite and Postgres in Spring 1989.



**Figure 6.** RAID-I mechanical configuration for magnetic disks. (The CPU is in a separate cabinet.) Note that disks can be placed in the back of the cabinet as well as the front, so the maximum capacity for this organization is 56 5.25" drives per rack. Using 340 MB Wren-IVs, the storage capacity of a rack is 19 gigabytes.

## 8. Discussion

We have discussed the RAID ideas several times with several groups of people, and some questions come up so commonly that we address them here.

### *1. How realistic are the assumptions of the independence of disk failures and constant failure rates? Is it a triumph of hope over experience?*

We wish we knew. There are no published papers showing the real lifetime failure rates of disks. Disk manufacturers estimate MTBF by accelerated life-cycle testing, assuming that failures are independent and the exponential failure model. If the reader knows of a source of such information, please contact the authors. It matters not who the manufacturer is, we just need empirical data on disks to evaluate the viability of alternatives for redundancy.

One problem that several magnetic disk manufacturers have mentioned is what we would call the "Pinto Effect;" a mistake is made in manufacturing process that is so disastrous that the disk manufacturer will recall all affected disks and replace them. The common theme is that the mistake is uncovered after the disks have been in the field for several months and the disks all fail within a short time of one another. One example was a manufacturer who glued together the two halves of an head-disk assembly, with this glue dissolving after the disks had been in the field for 18 months. Another example was that a new bactericide used in an air filter interacted with the disk surface so that many failures occurred six months later. A common cause of the Pinto Effect is that a supplier will change some component as a cost cutting measure without notifying the disk manufacturer, and disastrous consequences occur due to unforeseen interactions.

Although we desperately need real data on disk failures, we are performing studies of using models to estimate the impact of the Pinto Effect on RAID reliability.

### *2. What are the implications of having or not having a hot standby spare?*

A standby spare is an unused but electrically connected disk that can replace a failed disk in the system without human intervention. The major advantage of standby spare is reducing the mean time to repair (MTTR), with the disadvantage of increasing the cost and complexity of the system. As soon a disk fails, the system can immediately reconstruct the information onto the spare. Depending on the load on the other disks in the system and the capacity of the failed disk, the MTTR could be 10 to 60 minutes with a standby spare.

Without a standby spare, the MTTR will be significantly longer:

- The repair man must be contacted to bring a disk for replacement; this could be anywhere from 1 hour to 12 hours depending on the level of service.
- Disks have a limit as to how fast they can adjust to temperature changes; a typical specification is 5 degrees per hour. If the field engineer stores replacement disks in his car, then he may have to wait for the disk to acclimatize before installing it. This could be 0 to 4 hours. If disks are stored in the computer room to overcome acclimatization delays, this removes the cost advantage of not having standby spares.
- Even if the system is designed to allow hot spare insertion, in practice people responsible for a computer will not want someone to open cabinets and replace equipment on a working computer while many important jobs are running on the system. (A failed disk does not disable the system since the data can be reconstructed.) Thus sociological implications of manual replacement may extend MTTR to 24 to 72 hours.

### *3. If RAID turns out not to be the answer, what is?*

The slow seek and rotation delays of mechanical devices can easily be overcome with solid state memory. The recently introduced "Solid state disks" (SSD) provide more than ample speed to match the growth of CPUs, because mechanical devices are essentially removed from the hierarchy. The problem is that cost is essentially a factor of 10 to 20 times larger for solid state disks. In applications where a few hot spots dominate disk accesses, a more economic solution would be to automatically migrate data between SSD and magnetic disks to achieve higher performance at lower cost.

The cost of SSD over RAM is a battery. The advantage of a separate box with solid state memory over just larger main memory is that the SSD is as reliable as magnetic disks in the presence of bugs in the operating system and database software, while this is not the case for main memory. SSD also have the advantage that they are more easily multiported than main

memory, allowing other CPUs to access them in case of a CPU failure.

There is also a way to improve reliability without redundant data storage. Recent work has suggested that disk failures can be predicted. For decades field engineers have had diagnostics programs that they run to exercise a disk during a preventative maintenance cycle. The purpose is to decide whether or not the disk should be replaced even though it hasn't failed yet. This is clearly a failure prediction scheme. Lin and Siewiorek studied messages printed on the system console and were able to see indications of disk failures up to two weeks before the failure [Lin 1986]. DEC has a software product, called VAXSimPLUS, that takes advantage of the predictive nature of some failures. It purports to predict 90% disk failures far enough in advance to copy the data from the suspect disk onto a spare disk preventing any disruption of the data.

To achieve improvements in reliability similar to RAID's would take much higher accuracy of prediction than 90%. It is also clear that disks and interfaces that could provide an early warning system would be ideal parts for even less expensive RAID's.

#### 4. *Can a disk manufacturer successfully market a RAID?*

Disk manufacturers must live with the interfaces provided by computer systems; for example, SCSI interfaces, HBA interfaces, and operating systems. No matter how large a file is, UNIX will ask for it as a series of small (8 KB) blocks of data. Systems houses that can change the interfaces and the operating system have a much better opportunity to take advantage of the potential of RAID.

#### 5. *What are implications of RAID for standard interfaces like SCSI and IPI?*

RAID's will be constructed with hundreds of disks, increasing the physical distance between the computer and the furthest disk and the number of disks that must be attached to a common bus. Ideally future standards would allow longer distances for the connections, much higher transfer rates, and more devices per connection.

### 9. Conclusion

RAID's offer a cost effective option to meet the challenge of exponential growth in the processor and memory speeds. We believe the size reduction of personal computer disks is the key to the success of disk arrays, just as Gordon Bell argues that the size reduction of microprocessors is a key to the success in multiprocessors [Bell 85]. In both cases the smaller size simplifies the interconnection of the many components as well as packaging and cabling. While large arrays of mainframe processors are possible, it is certainly easier to construct an array from the same number of microprocessors (or PC drives). Just as Bell coined the term "multi" to distinguish a multiprocessor made from microprocessors, we use the term "RAID" to identify a redundant disk array made from personal computer disks.

With advantages in cost-performance, reliability, power consumption, and floor space, we expect RAID's to replace large drives in future I/O systems. There are, however, several open issues that may bear on the practicality of RAID's:

- *What will be the real lifetime of a RAID vs. MTBF calculated using the independent, exponential failure model?*
- *Will disk controller design limit RAID performance?*
- *How should 100 to 1000 disks be constructed and physically connected to the processor?*

### Acknowledgements

This work was supported by the National Science Foundation under grant # MIP-8715235 and the California MICRO program. We would like to thank the support of Sun Microsystems Incorporated, Impress/CDC, and our other industrial partners for the additional support of this research. Peter Chen was supported in part by an ONR fellowship and Garth Gibson was supported in part by both an IBM fellowship and a Computer Measurement Group award.

## References

- [Adaptec 87] AIC-6250, *IC Product Guide*, Adaptec, stock # DB0003-00 rev. B, 1987, p. 46.
- [Bell 85] Bell, C.G., "Multis: a new class of multiprocessor computers," *Science*, vol. 228 (April 26, 1985) 462-467.
- [Bitton 88] D. Bitton and J. Gray, "Disk Shadowing," *in press*, 1988.
- [Chen 89] P. Chen, "An Evaluation of Redundant Arrays of Disks using an Amdahl 5890," M.S. Report, 1989 (in preparation).
- [Dougkis 89] F. Dougkis and J. Ousterhout, "A log structured file system," Spring COMPCON 89, March 1, 1989, San Francisco, CA, (*in this report*).
- [Fujitsu 87] "M2361A Mini-Disk Drive Engineering Specifications," (revised) Feb., 1987, B03P-4825-0001A.
- [Gibson 89] G. Gibson, L. Hellerstein, R. Karp, R. Katz, and D. Patterson, "Error Correction in Large Disk Arrays," ASPLOS III, April 1989, Boston, MA.
- [Kim 86] M.Y. Kim, "Synchronized disk interleaving," *IEEE Trans. on Computers*, vol. C-35, no. 11, Nov. 1986.
- [Kurzweil 88] F. Kurzweil, "Small Disk Arrays - The Emerging Approach to High Performance," presentation at Spring COMPCON 88, March 1, 1988, San Francisco, CA.
- [Lin 86] T-T.Y. Lin and D.P. Siewiorek, "Architectural Issues for On Line Diagnostics in a Distributed Environment," *International Conference on Computer Design*, IEEE Computer Society, Rye Town, NY, October 1986.
- [Livny 87] M. Livny, S. Khoshafian, H. Boral, "Multi-disk management algorithms," *Proc. of ACM SIGMETRICS*, May 1987.
- [Park 86] A. Park and K. Balasubramanian, "Providing Fault Tolerance in Parallel Secondary Storage Systems," Department of Computer Science, Princeton University, CS-TR-057-86, Nov. 7, 1986.
- [Patterson 88] D. Patterson, G. Gibson, and R. Katz, "A Case for Redundant Arrays of Inexpensive Disks (RAID)," ACM SIGMOD conference proceedings, Chicago, IL., June 1-3, 1988, pp. 109-116. (Also appeared as Technical Report UCB/CSD 87/391, December 1987.)
- [Salem 86] K. Salem and Garcia-Molina, H., "Disk Striping," *IEEE 1986 Int. Conf. on Data Engineering*, 1986.
- [Schulze 89] M. Schulze, G. Gibson, R. Katz, and D. Patterson, "How Reliable is a RAID?," Spring COMPCON 89, March 1, 1988, San Francisco, CA, (*next paper in this report*).
- [Stonebraker 88] M. Stonebraker, R. Katz, D. Patterson, and J. Ousterhout, "The Design of XPRS", *Very Large Data Base Conference Proceedings*, August 1988, Long Beach, CA., pp. 318-330.

# HOW RELIABLE IS A RAID?<sup>1</sup>

*Martin Schulze<sup>2</sup>, Garth Gibson, Randy Katz, David Patterson*

Computer Science Division  
Electrical Engineering and Computer Science Department  
University of California, Berkeley  
Berkeley, CA 94720

**ABSTRACT:** Disk arrays provide the promise of greatly increased transfer bandwidth at low cost. But without additional data redundancy, an array can suffer from significantly degraded reliability. In this paper we more closely examine the reliability of RAID systems and find that although phenomenal reliability disk arrays can not be attained with data redundancy alone, RAID system reliability can be made better than conventional large disks with little extra hardware.

## 1. Introduction

A critical challenge for future computer systems is to match the predicted advances in processor performance with comparable improvements in the I/O system. Arrays of disk drives hold forth the promise of such improvements, by dramatically improving transfer bandwidth through the parallelism inherent in many disk arms with data spread across multiple drives. Small format drives are particularly attractive candidates for building disk arrays, because of their low cost and their high volumetric (MB/ft<sup>3</sup>) and power (MB/watt) efficiencies [VASU88].

The one significant drawback to using large numbers of disks to improve I/O performance is the impact it has on data reliability. More disks mean more disk failures and an increase in the probability of data loss. This observation is captured in the simple equation that relates the Mean Time To Failure (MTTF) of an array to the MTTFs of its component disks, assuming independent failures and a constant rate of failure:

$$MTTF_{disk\ array} = \frac{MTTF_{single\ disk}}{Number\ of\ Disks\ in\ the\ Array}$$

A disk array built from 49 small disks (ample capacity to replace conventional disks), each with an MTTF of

---

<sup>1</sup>Research Supported by National Science Foundation Grant # MIP-8715235 and the California MICRO Program with matching industrial support from SUN Microsystems, Inc. This paper will appear in COMPCON Spring 1989, San Francisco, CA, February 1989.

<sup>2</sup>Current Address: Digital Equipment Corporation, Colorado Springs, CO.

40,000 hours, would yield an overall MTTF of only 816 hours. This is significantly worse than the MTTF of a conventional mainframe disk with the same capacity. These conventional large disks have MTTFs of about 50,000 hours [BALA88].

Thus, arrays must be *redundant*, i.e., some capacity and bandwidth are sacrificed in order to redundantly store data, so that it is possible to reconstruct data lost in a failure. We have called such arrays **RAIDs**, for Redundant Arrays of Inexpensive Disks. This approach does not make array components more reliable, but, by making arrays tolerate component failures, it provides data availability comparable to conventional large disks. We view data availability as the reliability of the data. For simplicity, we call this the reliability of the array and assuming a constant failure rate, measure it with  $MTTF_{RAID}$ .

A RAID array is broken into a number of reliability *groups*, each containing extra "check" disks containing redundant data. This redundancy is intended to allow RAID systems to obtain a level of reliability at least equivalent to the conventional disk systems they are to replace. The following MTTF equation for a single error correcting RAID (i.e., an array that can tolerate a single failed disk in any reliability group, but not two failed disks) is reproduced from [PATT88]. It describes the reliability of a disk array (mean time to a failure that results in data loss) of  $n_G$  groups each with  $G$  data disks and 1 check disk when only failures in the disk devices are considered.

$$MTTF_{RAID} = \frac{(MTTF_{disk})^2}{n_G G (G+1) MTTR_{disk}}$$

In this equation  $MTTR_{disk}$  is the disk Mean Time To Repair and disk failure rates are assumed to be constant (i.e., disk lifetime is exponentially distributed). Following our earlier example of 49 data disks, this equation predicts that an array of 7 groups each with 7 data disks that have a 40,000 hour  $MTTF_{disk}$  and a 2 hour  $MTTR_{disk}$  will have an overall  $MTTF_{RAID}$  of 2,040,816 hours. Fast repair is facilitated by onsite spare (inexpensive) disks and operating system directed online replacement and reconstruction. Unfortunately, this phenomenal  $MTTF_{RAID}$ , 239 years, is a very optimistic estimate of reliability because failures of disk support hardware have been ignored.



The purpose of this paper is to develop a better understanding of RAID reliability than the simple MTTF calculation just introduced. The reliability of other components, such as power supplies, controller electronics, cables, and fans, should also be considered for they affect the overall reliability of the array. Our goal is not “non-stop” reliability, but rather a level of reliability for large numbers of disks that is comparable to that of a single conventional disk. However, if failure rates of disks and support hardware are constant, we shall see that seemingly excessively large MTTFs may be desirable to secure a low probability of data loss.

The rest of this paper is organized as follows. In the next section we briefly review basic reliability definitions and examine sources of non-catastrophic disk failures. We next examine the exponential failure model and conclude that users’ perceptions of MTTFs are not really what MTTFs mean. We then develop a model of reliability for array systems that includes support hardware. This leads us to hardware reliability groups, which contain redundant hardware components to further improve the array reliability.

## 2. How Disks Fail

We adopt definitions from the fault-tolerant research community [MAXI88]. A *failure* is a detectable physical change to hardware. Failures may be repaired by the replacement of a physical component. A *fault* is an event which interferes with normal operation and can be either *soft (transient)*, i.e., not readily repeatable, or *hard*, i.e., repeatable with high probability. Hard faults may be caused by failures, while soft faults are more likely caused by environmental factors or insufficient design margins. An *error* is a manifestation of a fault by an incorrect value. Errors, therefore, can be either soft or hard.<sup>3</sup> In this paper we are mainly concerned with catastrophic failures; failures that render a device module inoperable (such as head crashes or read/write electronics failures), but we begin by exonerating non-catastrophic failures and faults.

Disk drive manufacturers have identified a few types of errors associated with the servo system (positioning the heads) and the read/write system as critical to customer satisfaction. Typical specifications for the

---

<sup>3</sup> In some systems recoverable errors are called *soft* and unrecoverable errors are called *hard*. This differs from our usage mainly through repeatable errors that are recoverable by error correcting codes.

Table I -- Disk Drive Error Rates			
Type of Error	Average Error Rate	Recovery	Consequences
Recoverable Seek Error	<1 error in $10^6$ seeks	retry	none
Transient Recoverable Data Error	<1 error in $10^{10}$ bits read	retry or ECC	none
Repeatable Recoverable Data Error	< 1 error in $10^{12}$ bits read	ECC	Data rewritten to relocated sector.
Unrecoverable Data Error	< 1 error in $10^{14}$ bits read	none	One sector's data lost. Sector marked bad.
Miscorrected Data Error	< 1 error in $10^{21}$ bits read	none	One sector's data incorrectly read.

occurrence rates of these types of errors are shown in Table I [CDC 88, QUAN87].

A recoverable seek error is a seek in which the drive does not locate the desired cylinder on the first try, but is successful during a retry (if it is never successful then a catastrophic failure has occurred). A data error is defined as one sector read incorrectly, as detected by an Error Correcting Code (ECC). Random recoverable data errors are soft errors usually related to the signal-to-noise ratio of the system. Repeatable recoverable errors are hard errors, most often due to media defects, that can be corrected by ECC. Unrecoverable data errors lose data because the sector is detectably too damaged to recover by ECC. Miscorrected data errors occur when ECC was incorrectly not invoked or has resulted in incorrect data.

To get a feeling for the magnitude of these error rates, consider a disk that is performing 50 seeks/sec and reading 512KB/sec sustained. The mean time to next failure for each of these error types is: for recoverable seek errors, 5.6 hours; for random recoverable data errors, 40 minutes; for repeatable recoverable data errors, 2.8 days; for unrecoverable data errors, 276 days; and for miscorrected data errors, 7.6 million years.

The first three of these are recoverable without user intervention and the last is negligible when it is compared to each disk's 5 year mean time to catastrophic failure. Only unrecoverable data errors appear to pose a problem, but these can be dealt with in the same manner as a catastrophic drive failure – the sector can be reconstructed from redundant data in the array. We now turn our attention to just how low we should try to make a disk array's catastrophic failure MTTF.

### 3. Exponential Failure Model

Many disk drives offer catastrophic failure reliabilities specified as MTTF = 30,000 to 50,000 hours of normal usage. Disk users tend to interpret these specifications as implying that their disks will fail after this many hours of operation. Unfortunately, this is not the case. If disk lifetimes are truly exponential, i.e., failure rates are constant, with mean lifetime equal to the MTTF, then there is a large probability that the disk will fail before the MTTF is reached. Table II summarizes this observation. The actual formula is:

$$\text{Prob(exactly } k \text{ failures in MTTF)} = \frac{1}{(e^k k!)}$$

where  $e$  is the base of the natural logarithm.

This means that some customers will have one or more failures within a small fraction of the quoted MTTF. To avoid customer distress, disk manufacturers may underrepresent their MTTFs. Alternatively, in

Table II -- Exponential Lifetimes compared to MTTF			
Percent of All Disks	# of Failures Experienced Within MTTF	Cumulative Percent	# of Failures Experienced Within MTTF
36.8	0	100.0	0 or more
36.8	1	63.2	1 or more
18.4	2	26.4	2 or more
6.1	3	8.0	3 or more
1.5	4	1.9	4 or more
0.3	5	0.4	5 or more
0.05	6	0.06	6 or more
0.007	7	0.008	7 or more

building a disk system, the goal should be to make the probability of data loss within a reasonable interval as low as possible. Thus we should strive for MTTFs substantially higher than the expected useful product lifetime. However, the overall reliability of users' computation also depends on main memory and CPU failure rates and even more so on software quality [GRAY85], so we should not be too extravagant with the design of an IO subsystem.

#### 4. RAID Reliability Revisited

The redundancy of RAID is an application of a fault tolerant technique to address the problem of data loss due to disk drive failures. Small disk drives are not standalone units, but require support hardware: power supplies, SCSI (Small Computer System Interface) Host Bus Adapters (HBAs), cooling equipment, and cabling. To get a more accurate picture of RAID reliability, all parts of the array system should be considered.

The Berkeley RAID is based on the concept of *parity group*, i.e., a group of disks sharing a common parity check disk. When this RAID is implemented with a shared interconnect such as SCSI, a second type of grouping emerges that is based on the *SCSI group*, i.e., a group of disks sharing a common SCSI cable and HBA. There is also the *power group*, a group of disks sharing a common power supply, and the *cooling group*, a group of disks sharing a common fan. The interaction of these groups has a major influence on overall RAID reliability.

Figure 1 shows the reliability of various components of a SCSI based RAID. Note that the RAID MTTF equation presented in Section 1 considered only the reliability of disks. To build a realistic system assembly would require eight SCSI HBAs with cables, eight 300 Watt power supplies (each with a power cable), and eight fans for cooling. The overall RAID reliability should take account of these additional components and their independent rates of failure. Of greatest concern is the external power grid; without battery backup your system is at the mercy of the power company and can expect an overall MTTF no better than 2 months! Assuming that the power supply has battery backup but that any failure in the support components may cause data loss (a pessimistic assumption),  $MTTF_{RAID}$  would be revised to 5734 hours. This is about

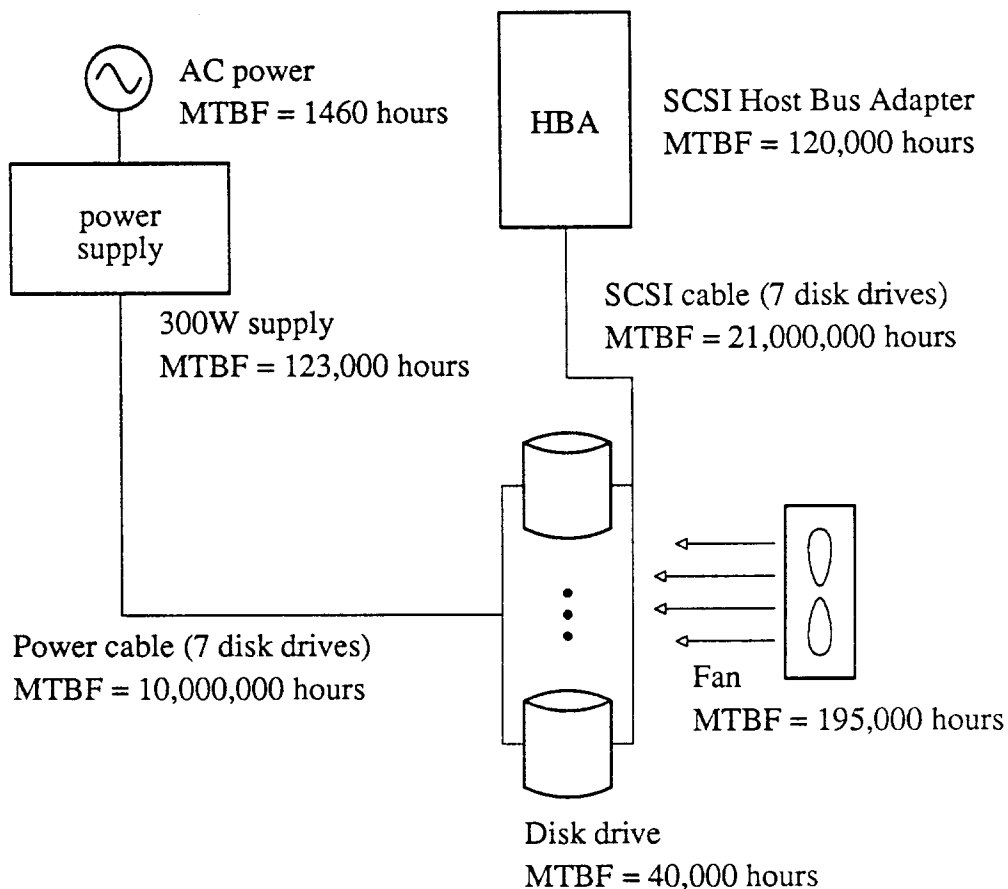


Figure 1: Reliability of RAID Components

A disk subsystem is more than just disks. This figure illustrates typical reliabilities for each part. The MTTFs for disk drives and Host Bus Adapters are estimates quoted from representative manufacturers [CDC88, Moren88]. MTTFs for fans are estimated from MIL-HDBK-217E part code 801 (electric motor, < 1 horsepower) with 4 solder connections. MTTFs for SCSI cables are estimated from MIL-HDBK-217E part code 1105 (printed wiring board connector) with 50 active pins and 50 milliamps per pin with 0.04 mate/unmate cycles per 1000 hours. Power supply MTTFs are from MIL-HDBK-217D [BARD86] and MTTFs for power cables are estimated from MIL-HDBK-217E part code 1103 (power connector) with 4 active pins and 2.5 amps per pin with 0.04 mate/unmate cycles per 1000 hours. The MTTF of the external power grid is taken from Gray [GRAY85].

239 days and represents a factor of 356 decrease in MTTF from the simple estimate that considered only disk drive failures!

However, by judicious placement of parity, SCSI, power, and cooling “groups”, we can do much better. If parity groups are mapped onto the disk array orthogonal to SCSI, power, and cooling groups, then no single hardware failure will cause data loss. We see this in Figure 2 because the loss of any complete

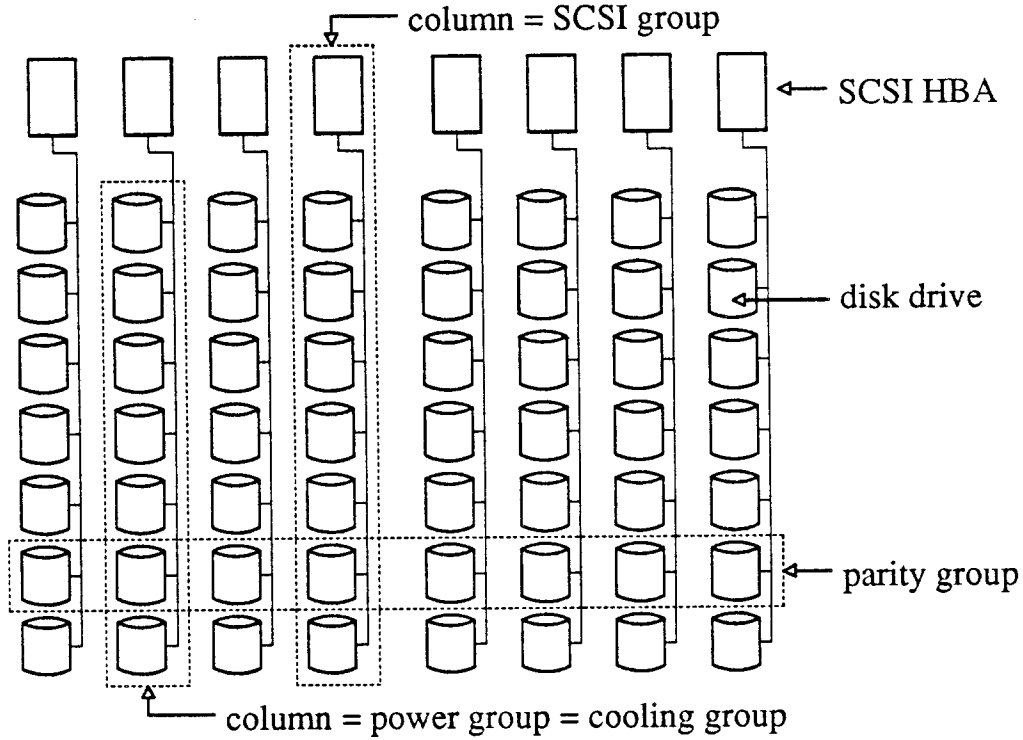


Figure 2: Groupings in a Disk Array

Data reliability groups are organized orthogonal to a SCSI string, while hardware reliability groups are organized around the string. This improves reliability by insuring that no single data or hardware failure will lead to data loss.

column amounts to the loss of a single disk in each parity group and each of these is recoverable. This scheme does not have any explicit fault tolerance of the support hardware, but uses the redundancy of parity groups to protect against support hardware failures as well as disk failures. For this to work, support hardware failures should be repaired with little or no interruption of service. Since these failures involve a variety of types of equipment that may not be easily replaced, a separate, and probably longer, mean time to repair,  $MTTR_{column}$ , is used.

For the RAID of Figure 2,  $MTTF_{RAID}$  can be estimated as

$$\frac{(MTTF_{disk})^2}{n_G G (G+1) MTTR_{disk} (1 + \alpha_F \frac{1 + \alpha_R}{\alpha_R} + \frac{\alpha_F^2}{n_G \alpha_R})}$$

$$\text{where } \alpha_F = \frac{MTTF_{disk}}{MTTF_{column}} \text{ and } \alpha_R = \frac{MTTR_{disk}}{MTTR_{column}}.$$

Note that this formula reduces to the original as  $MTTF_{column}$  goes towards infinity.

To illustrate this model, reconsider our earlier example. Assuming each column has one SCSI HBA, one 300 Watt power supply, one fan, one SCSI cable, and one power cable, then  $MTTF_{column}$  is 46,000 hours (calculated by summing the failure rates of these support components). If the average time to repair support hardware failures in a column,  $MTTR_{column}$ , is 72 hours, then  $MTTF_{RAID}$  is 55,000 hours. Although 55,000 hours is still a far cry from the simple estimate of 2,000,000 hours, it does meet our goal of exceeding the reliability of an individual conventional disk. In fact, conventional disk reliabilities also exclude interconnect and host bus adapter support hardware, so, using a SCSI interconnect and HBA, a conventional single disk

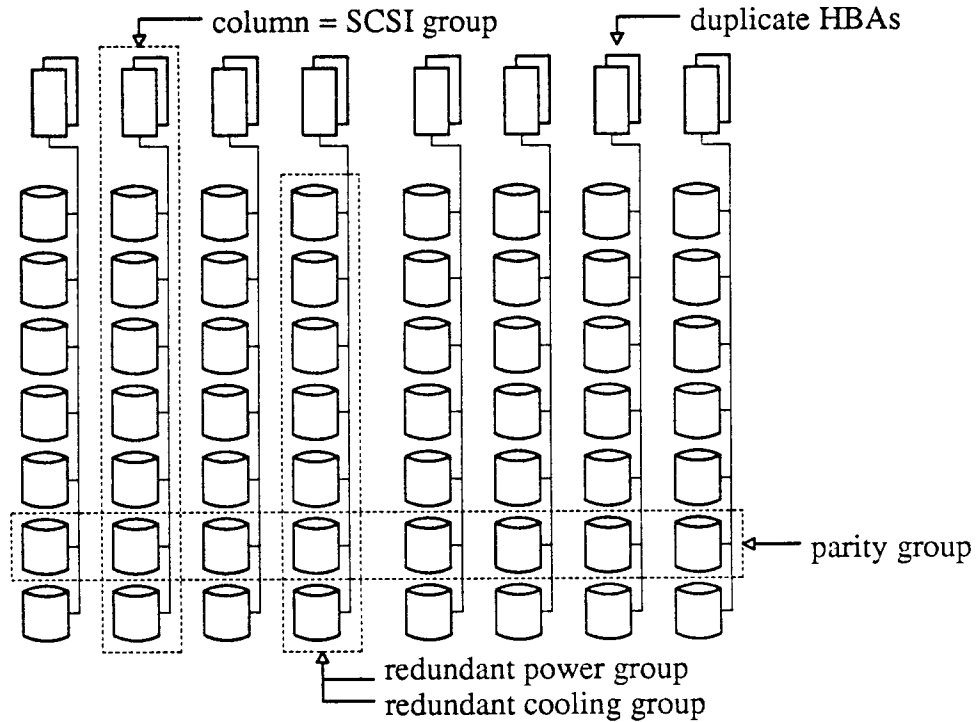


Figure 3: Redundant Groupings in a Disk Array

Adding redundant hardware to the column can increase its MTTF and thus the reliability of the whole array.

subsystem built around 50,000 hour MTTF disk will have an overall MTTF of 35,235 hours (large conventional disks have built in power and cooling).

Just as data redundancy dramatically improves the reliability of disk data, additional hardware redundancy should improve the reliability of the support hardware. For example, power and cooling groups can be made redundant, or SCSI HBAs can be duplicated on each SCSI cable (see Figure 3). By adding redundant hardware components,  $MTTF_{column}$  can be increased, thus further increasing  $MTTF_{RAID}$ . Table III shows the effects on RAID reliability of making various combinations of the support hardware redundant. We show both how far each combination is from the original  $MTTF_{RAID}$  estimate, as a percent, and how this compares to a conventional single disk subsystem reliability of 35,235 hours, also as a percent. We should stress here that our model pessimistically assumes that fan failures render an entire column unavailable and two fan failures close together cause the loss of data.

Table III -- RAID Data Reliability Considering Support Hardware

Configuration	$MTTF_{RAID}$	Percent of Maximum	Percent of Conventional
No Redundancy	816 hrs	0.04%	2.1%
Simple RAID	5,734 hrs	0.3%	16%
(Orthogonal) RAID	55,000 hrs	3%	156%
RAID + Redundant Fans	73,000 hrs	4%	208%
RAID + Redundant Power Supplies	90,000 hrs	4%	255%
RAID + Redundant HBAs	91,000 hrs	5%	259%
RAID + Redundant Power and Fans	144,000 hrs	7%	409%
RAID + Redundant HBAs and Fans	148,000 hrs	7%	418%
RAID + Redundant Power and HBAs	225,000 hrs	11%	639%
RAID + Redundant Power, HBAs, Fans	1,650,000 hrs	81%	4680%



## 5. Summary and Conclusions

RAID is a technique for increasing I/O bandwidth by spreading data across large numbers of small format disk drives organized as an array. Because a small number of conventional large format disks are replaced by much larger numbers of small disks, the reliability of the disk system is greatly reduced. To deal with the problem of substantially more frequent disk failures, RAID systems trade off some of the increased bandwidth and storage capacity for redundant data storage, so lost data can be reconstructed through parity calculations. An optimistic analysis, focusing on disk reliability alone, would indicate that RAID systems can be made very much more reliable than conventional systems at a very modest cost of extra "check" disks.

This paper has presented a more careful reliability analysis that clearly shows that the rest of the system components, such as the SCSI HBAs, power supplies, and fans cannot be ignored. We have presented a scheme in which the system support components are organized into groups orthogonal to the data redundancy groups, thus guaranteeing that no single disk OR component failure will permanently lose data. This approach yields disk arrays with reliability about 50% greater than conventional disk subsystems. If further reliability improvements are sought, various parts of the support hardware can be made redundant. We have shown how these affect the design of a 49 data disk RAID. From this example we see that redundancy in the relatively inexpensive fans and power supply can yield overall MTTFs of about 4 times conventional disk subsystems.

## 6. References

- [BALA88] Balanson, R., "GPD Products and Technology Directions," IBM Fellowship Conference Presentation, San Jose, CA, (November 1988).
- [BARD86] Bardos, P., "The Reliability of Switch Mode Power Supplies," *Electronic Engineering*, V. 58, N 715, (July 1986), pp. 37-44.
- [CDC 88] Control Data Corporation, "Product Specification for WREN IV SCSI Model 94171-344," Control Data OEM Product Sales, Minneapolis, MN, (January 1988).

- [GRAY85] Gray, J., "Why Do Computers Stop and What Can Be Done About It?," Tandem Technical Report 85.7, (June 1985).
- [MAXI88] Maxion, R.A., D.P. Siewiorek, "Symptom-directed diagnosis of distributed computing systems," 1986/1987 *Research Review*, Computer Science Department, Carnegie Mellon, 1988.
- [MIL 86] U. S. Department of Defense, *Military Handbook: Reliability Prediction of Electronic Equipment*, MIL-HDBK-217E, (October 1986).
- [MORE88] Moren, W. D., Ciprico, Inc., private communication, (July 1988).
- [PATT88] Patterson, D. A., G. Gibson, R. H. Katz, "A Case for Redundant Arrays of Inexpensive Disks (RAID)," ACM SIGMOD Conference, Chicago, IL, (June 1988).
- [SCHU88] Schulze, M. E., "Considerations in the Design of a RAID Prototype," M.S. Report, U. C. Berkeley Computer Science Division, (August 1988).
- [QUAN87] Quantum Corporation, "OEM/Programmers Manual for Q200 Series Disk Drives," Milpitas, CA, (May 1987).
- [VASU88] Vasudeva, A., "A Case for Disk Array Storage Systems," Proc. Systems Design and Networks Conference, Santa Clara, CA, (April 1988).