# Appendix D

Authors: John Hennessy & David Patterson
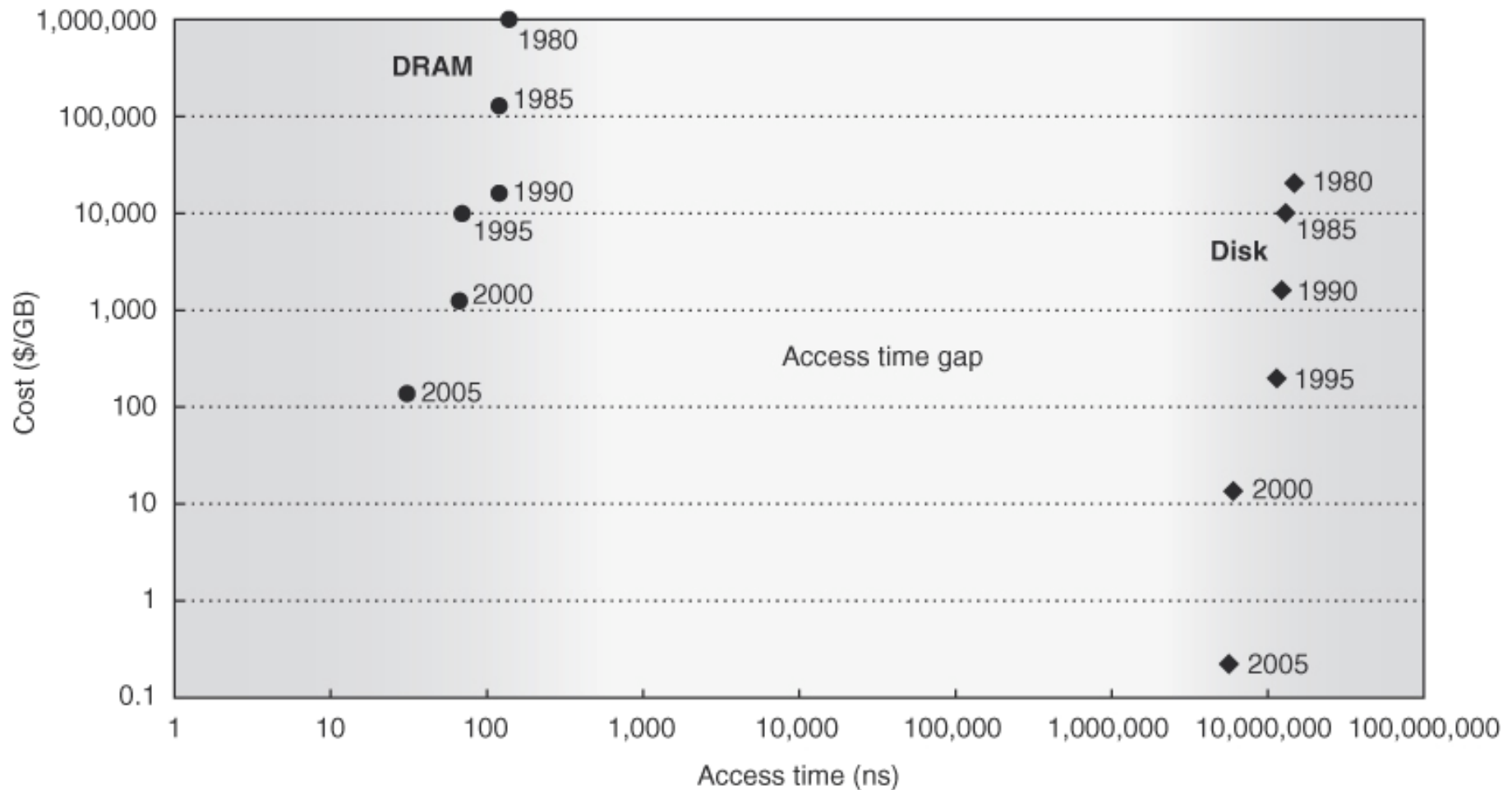
**Figure D.1 Cost versus access time for DRAM and magnetic disk in 1980, 1985, 1990, 1995, 2000, and 2005.** The two-order-of-magnitude gap in cost and five-order-of-magnitude gap in access times between semiconductor memory and rotating magnetic disks have inspired a host of competing technologies to try to fill them. So far, such attempts have been made obsolete before production by improvements in magnetic disks, DRAMs, or both. Note that between 1990 and 2005 the cost per gigabyte DRAM chips made less improvement, while disk cost made dramatic improvement.
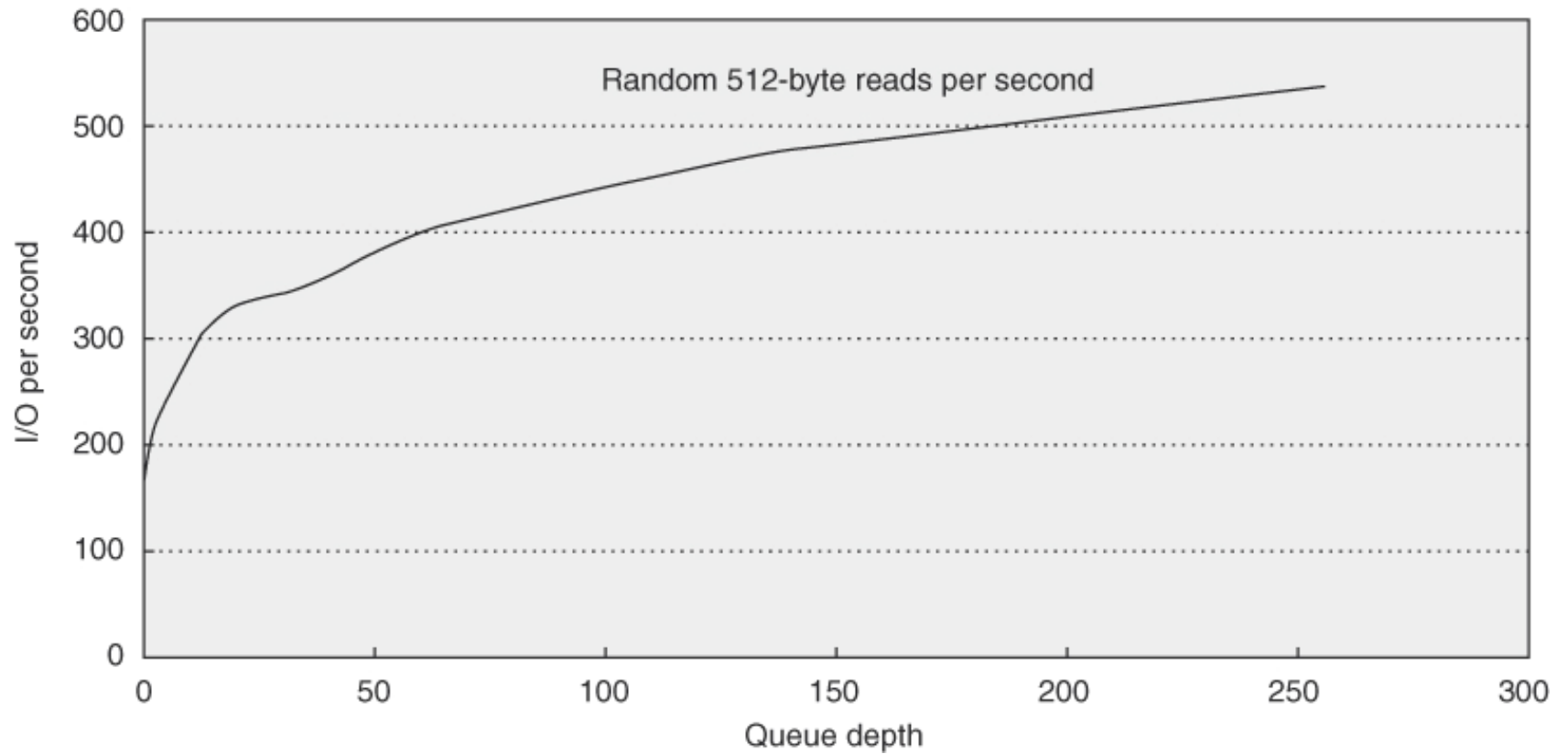
2

**Figure D.2 Throughput versus command queue depth using random 512-byte reads.** The disk performs 170 reads per second starting at no command queue and doubles performance at 50 and triples at 256 [Anderson 2003].

**Figure D.5 Row diagonal parity for $p = 5$, which pro-tects four data disks from double failures [Corbett et al. 2004].** This figure shows the diagonal groups for which parity is calculated and stored in the diagonal parity disk. Although this shows all the check data in separate disks for row parity and diagonal parity as in RAID 4, there is a rotated version of row-diagonal parity that is analogous to RAID 5. Parameter $p$ must be prime and greater than 2; however, you can make $p$ larger than the number of data disks by assuming that the missing disks have all zeros and the scheme still works. This trick makes it easy to add disks to an existing system. NetApp picks $p$ to be 257, which allows the system to grow to up to 256 data disks.
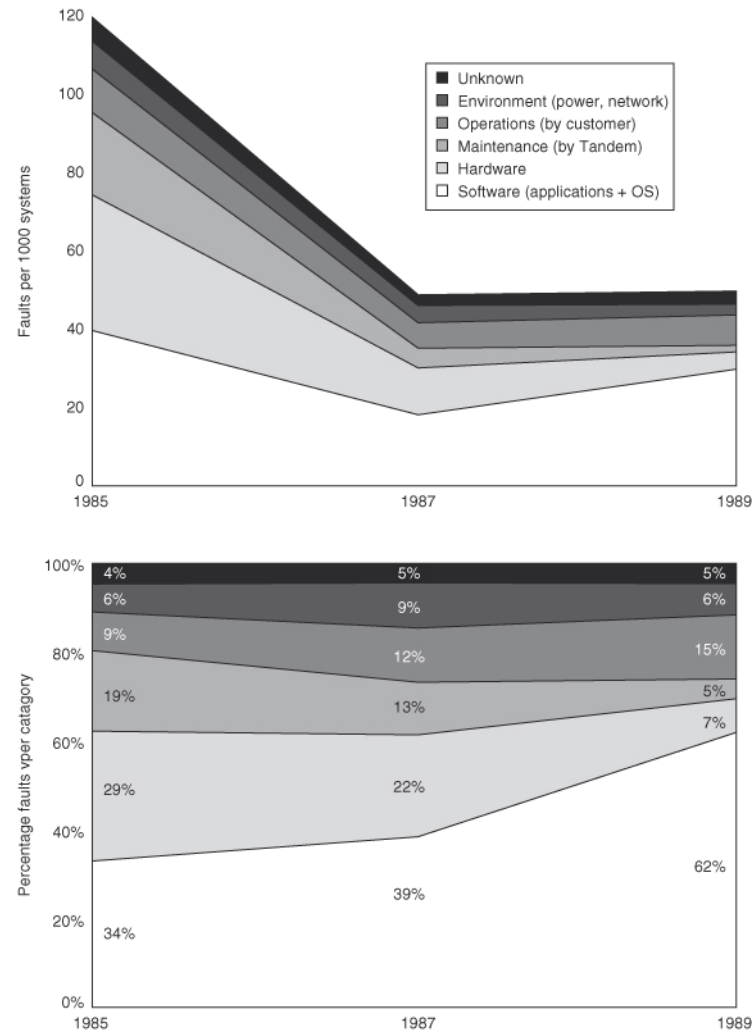
**Figure D.7 Faults in Tandem between 1985 and 1989.** Gray [1990] collected these data for fault-tolerant Tandem Computers based on reports of component failures by customers.
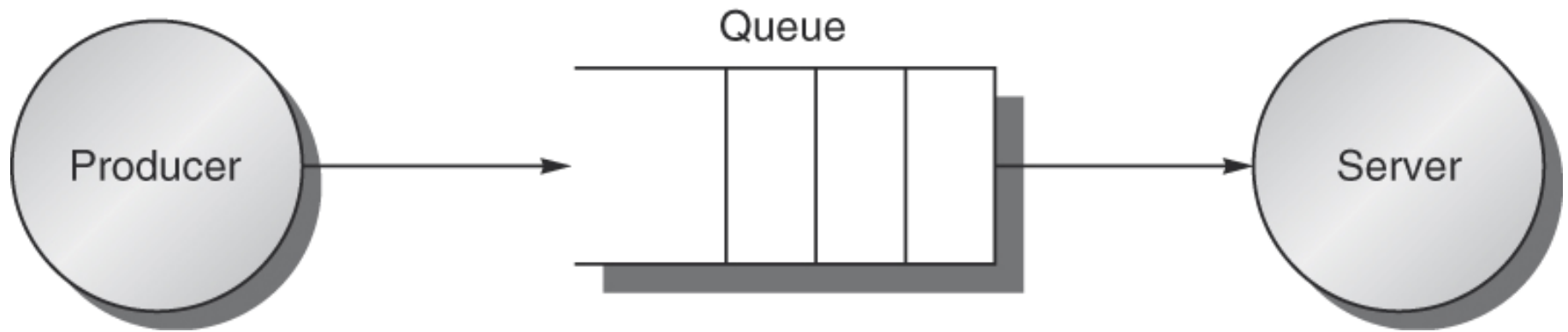
**Figure D.8 The traditional producer-server model of response time and throughput.** Response time begins when a task is placed in the buffer and ends when it is completed by the server. Throughput is the number of tasks completed by the server in unit time.
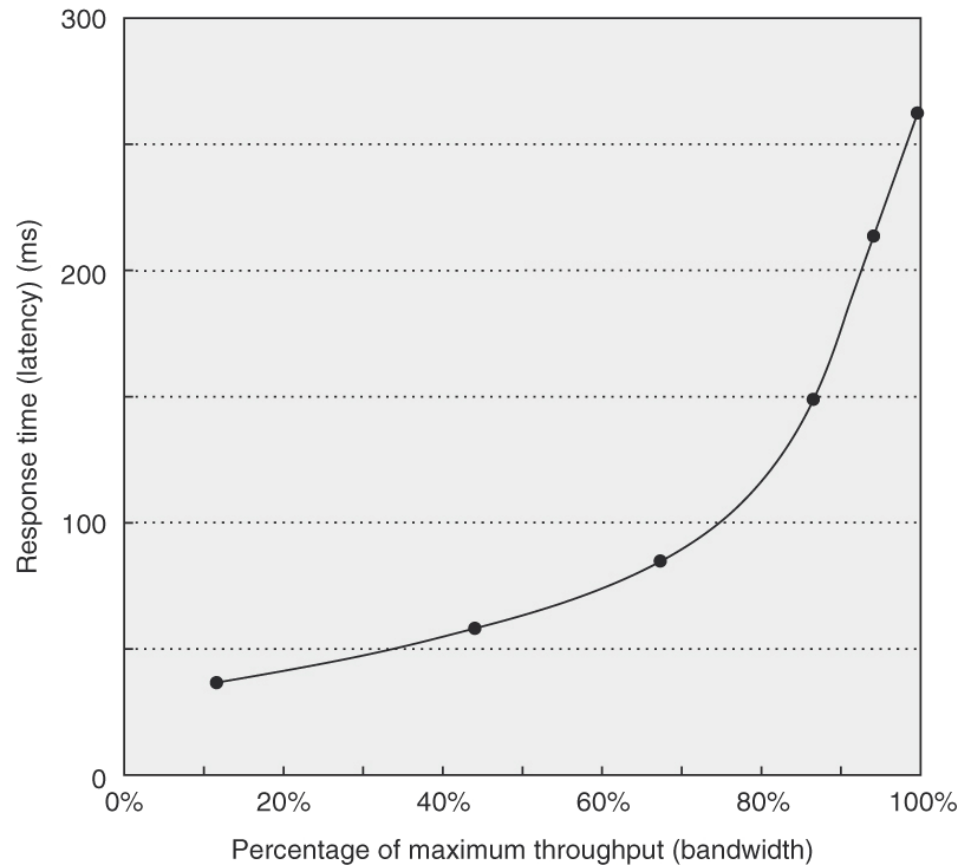
**Figure D.9 Throughput versus response time.** Latency is normally reported as response time. Note that the mini-mum response time achieves only 11% of the throughput, while the response time for 100% throughput takes seven times the minimum response time. Note also that the independent variable in this curve is implicit; to trace the curve, you typically vary load (concurrency). Chen et al. [1990] collected these data for an array of magnetic disks.
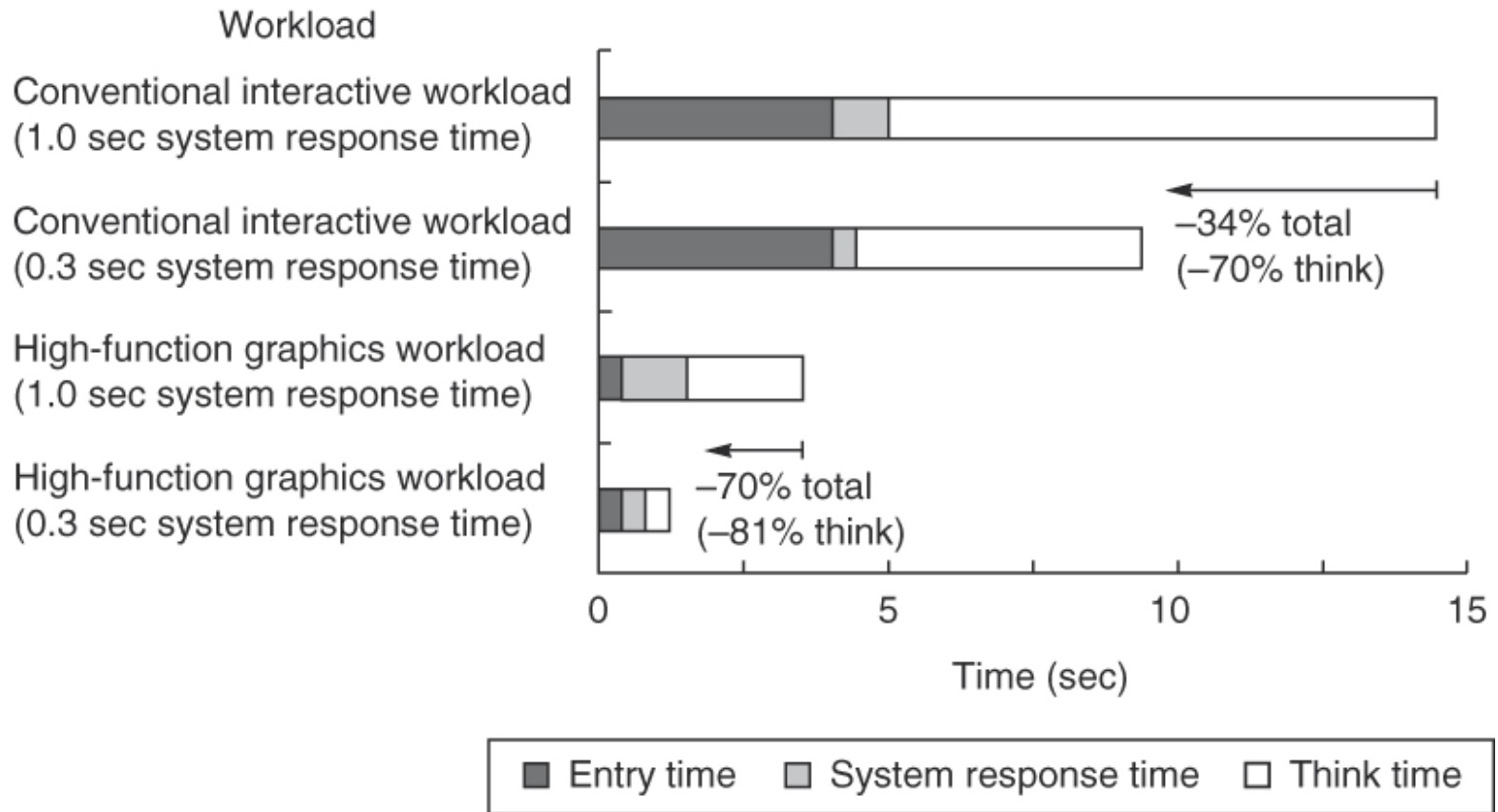
**Figure D.10 A user transaction with an interactive computer divided into entry time, system response time, and user think time for a conventional system and graphics system.** The entry times are the same, independent of system response time. The entry time was 4 seconds for the conventional system and 0.25 seconds for the graphics system. Reduction in response time actually decreases transaction time by more than just the response time reduction. (From Brady [1986].)
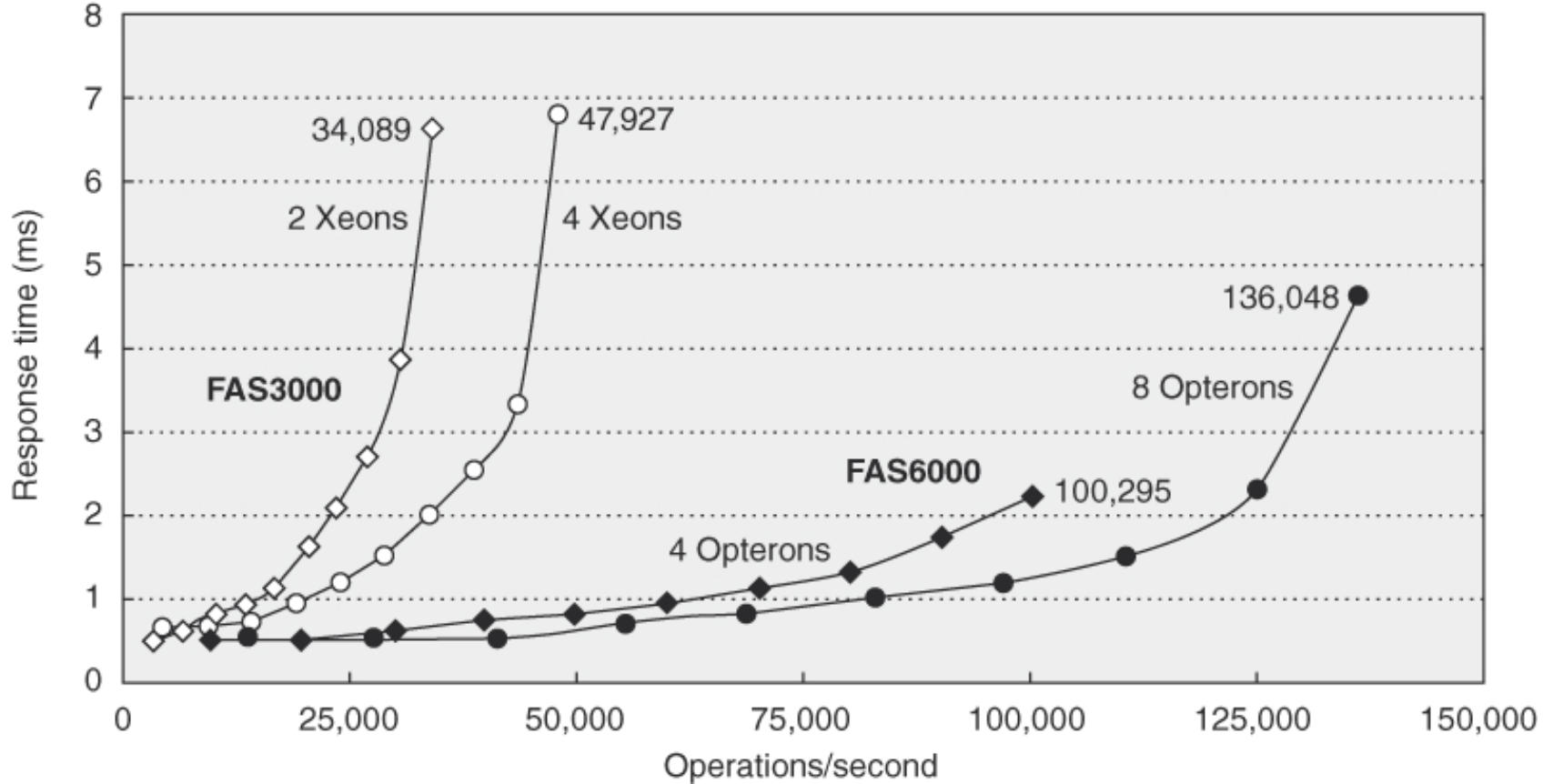
**Figure D.13 SPEC SFS97_R1 performance for the NetApp FAS3050c NFS servers in two configurations.** Two processors reached 34,089 operations per second and four processors did 47,927. Reported in May 2005, these systems used the Data ONTAP 7.0.1R1 operating system, 2.8 GHz Pentium Xeon microprocessors, 2 GB of DRAM per processor, 1 GB of nonvolatile memory per system, and 168 15K RPM, 72 GB, Fibre Channel disks. These disks were connected using two or four QLogic ISP-2322 FC disk controllers.
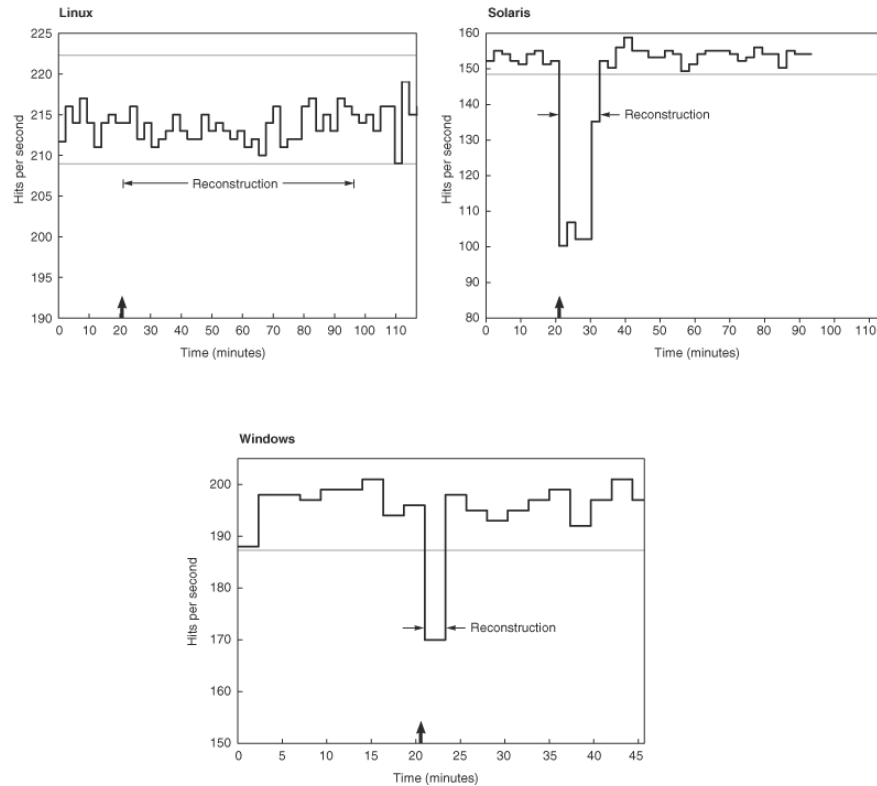
**Figure D.14 Availability benchmark for software RAID systems on the same computer running Red Hat 6.0 Linux, Solaris 7, and Windows 2000 operating systems.** Note the difference in philosophy on speed of reconstruction of Linux versus Windows and Solaris. The *y*-axis is behavior in hits per second running SPECWeb99. The arrow indicates time of fault insertion. The lines at the top give the 99% confidence interval of performance before the fault is inserted. A 99% confidence interval means that if the variable is outside of this range, the probability is only 1% that this value would appear.
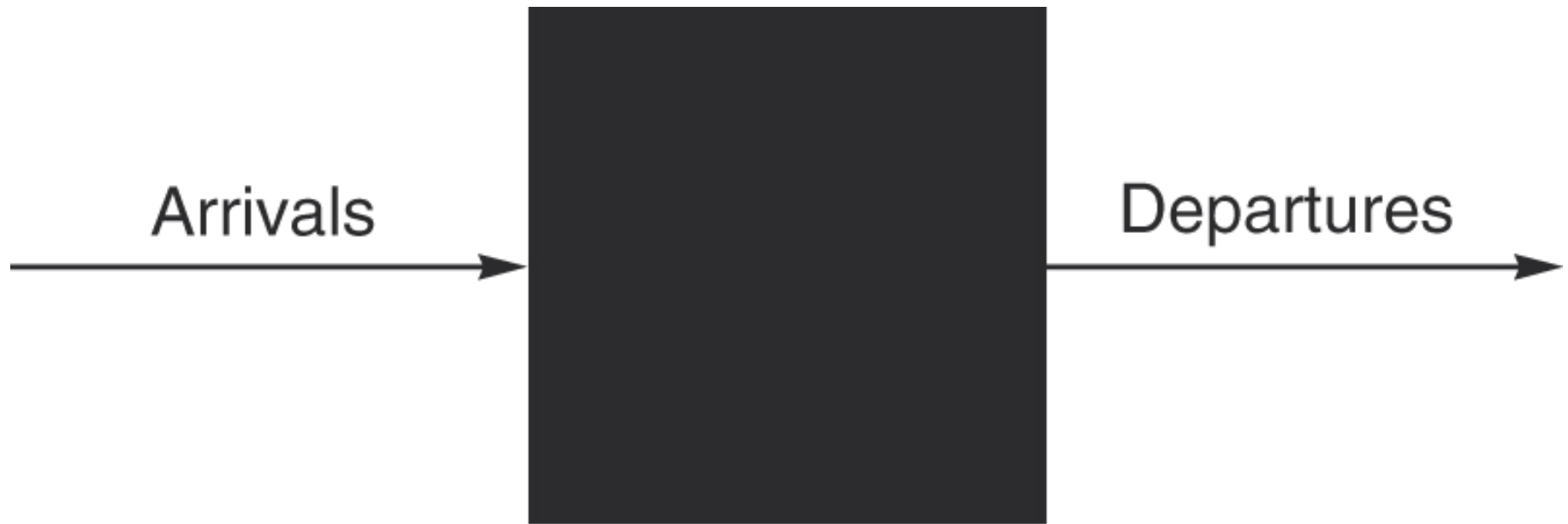
10

**Figure D.15 Treating the I/O system as a black box.** This leads to a simple but important observation: If the system is in steady state, then the number of tasks entering the system must equal the number of tasks leaving the system. This *flow*-balanced state is necessary but not sufficient for steady state. If the system has been observed or measured for a sufficiently long time and mean waiting times stabilize, then we say that the system has reached steady state.
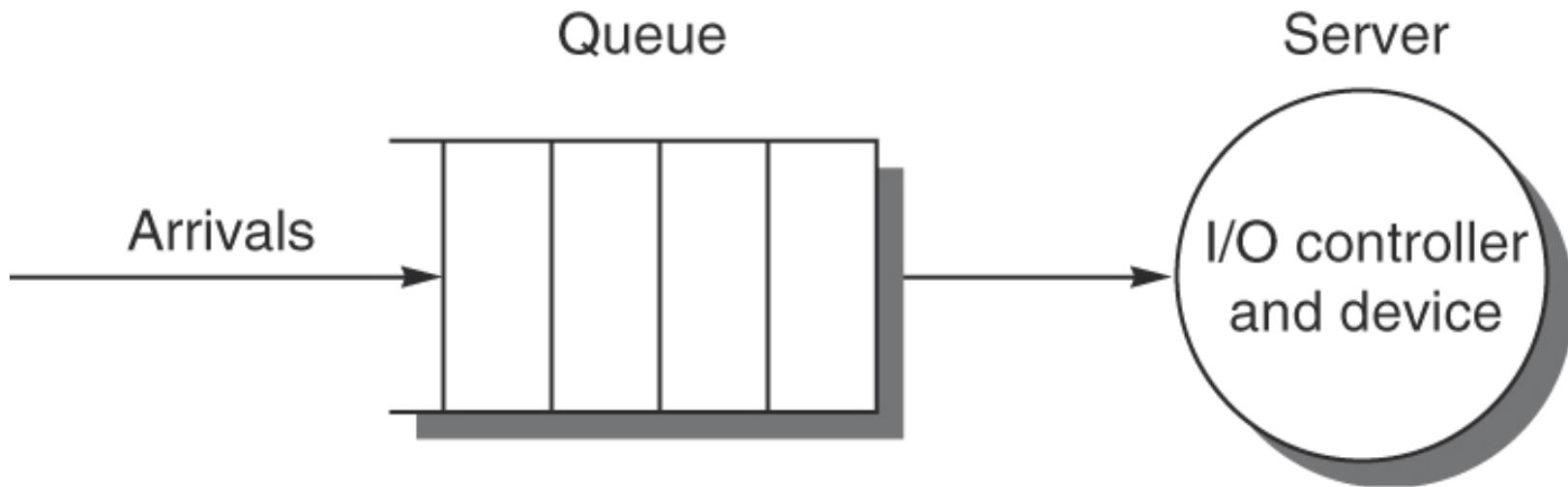
11

**Figure D.16 The single-server model for this section.** In this situation, an I/O request "departs" by being completed by the server.
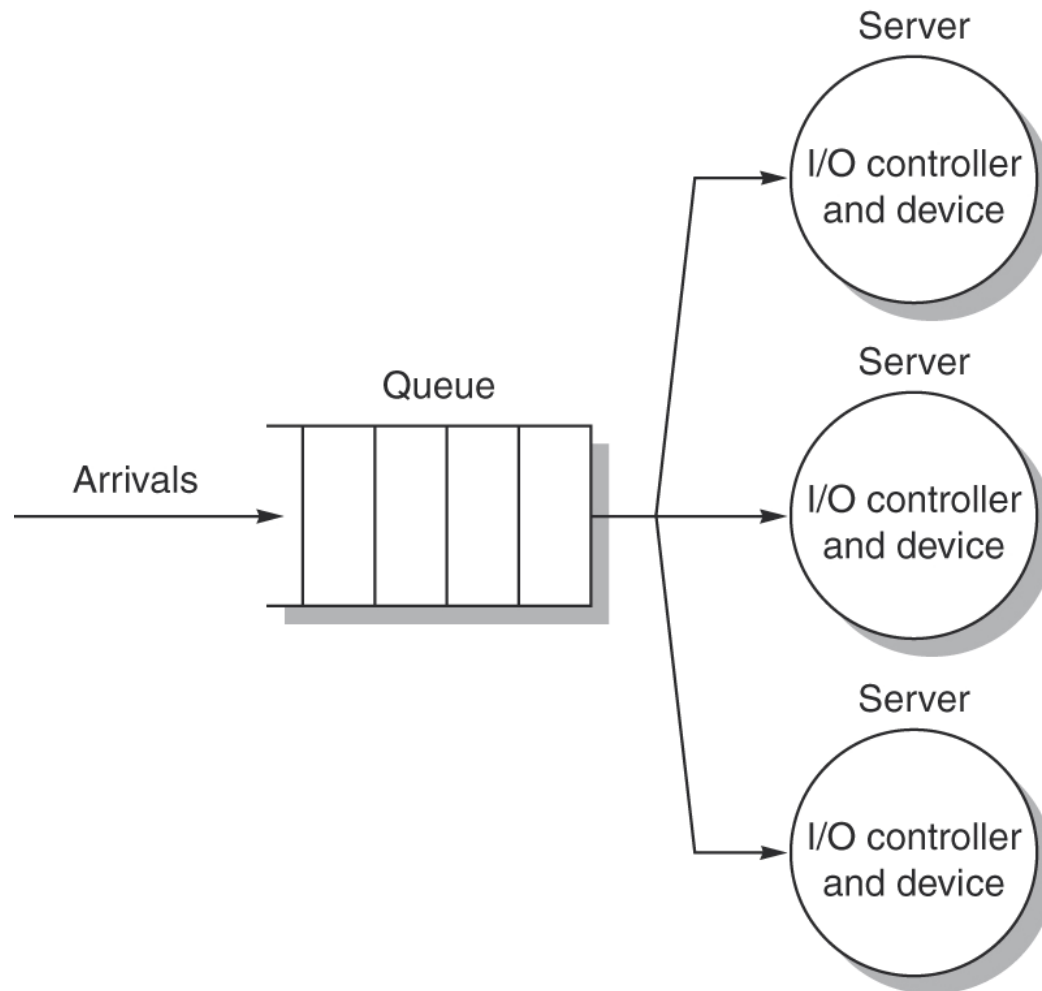
**Figure D.17 The M/M/m multiple-server model.**

**Figure D.19 The TB-80 VME rack from Capricorn Systems used by the Internet Archive.** All cables, switches, and displays are accessible from the front side, and the back side is used only for airflow. This allows two racks to be placed back-to-back, which reduces the floor space demands in machine rooms.
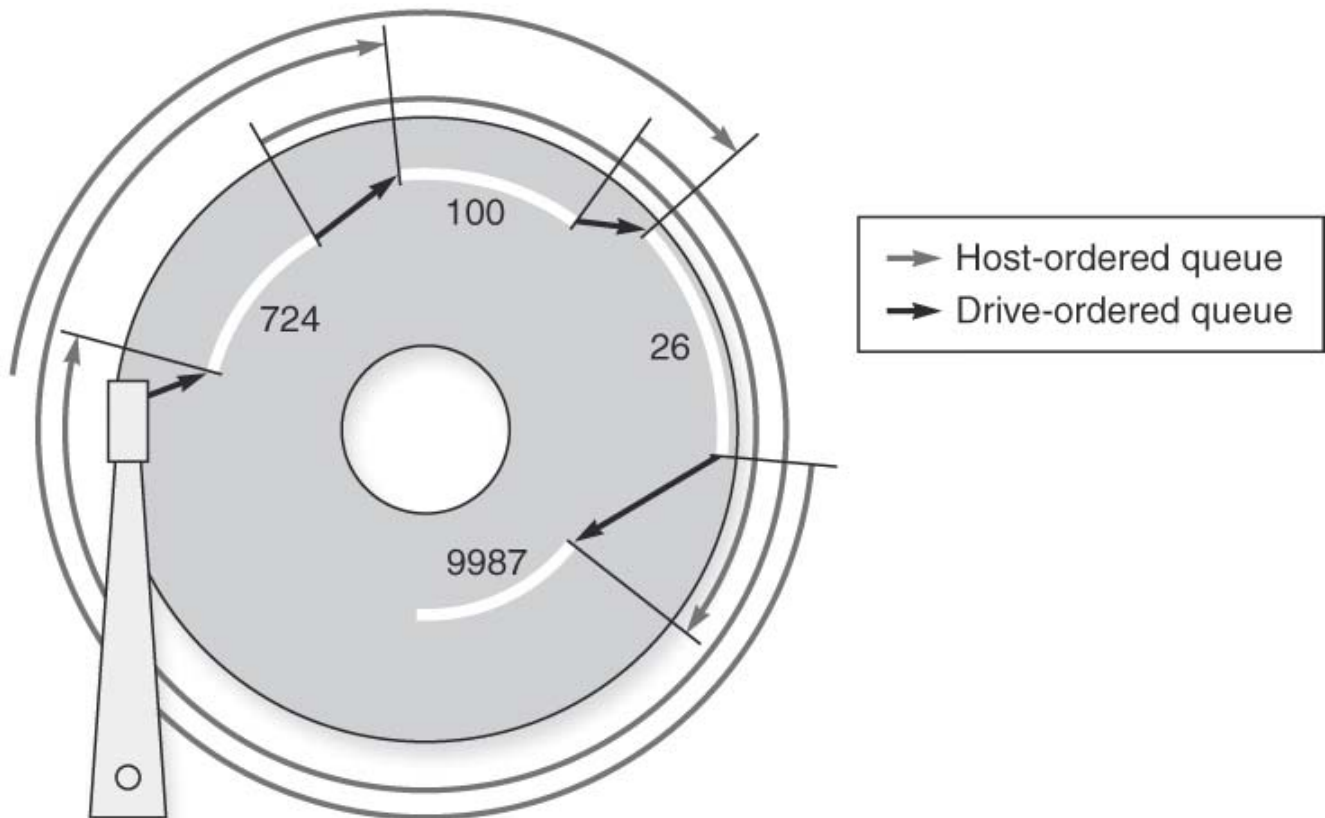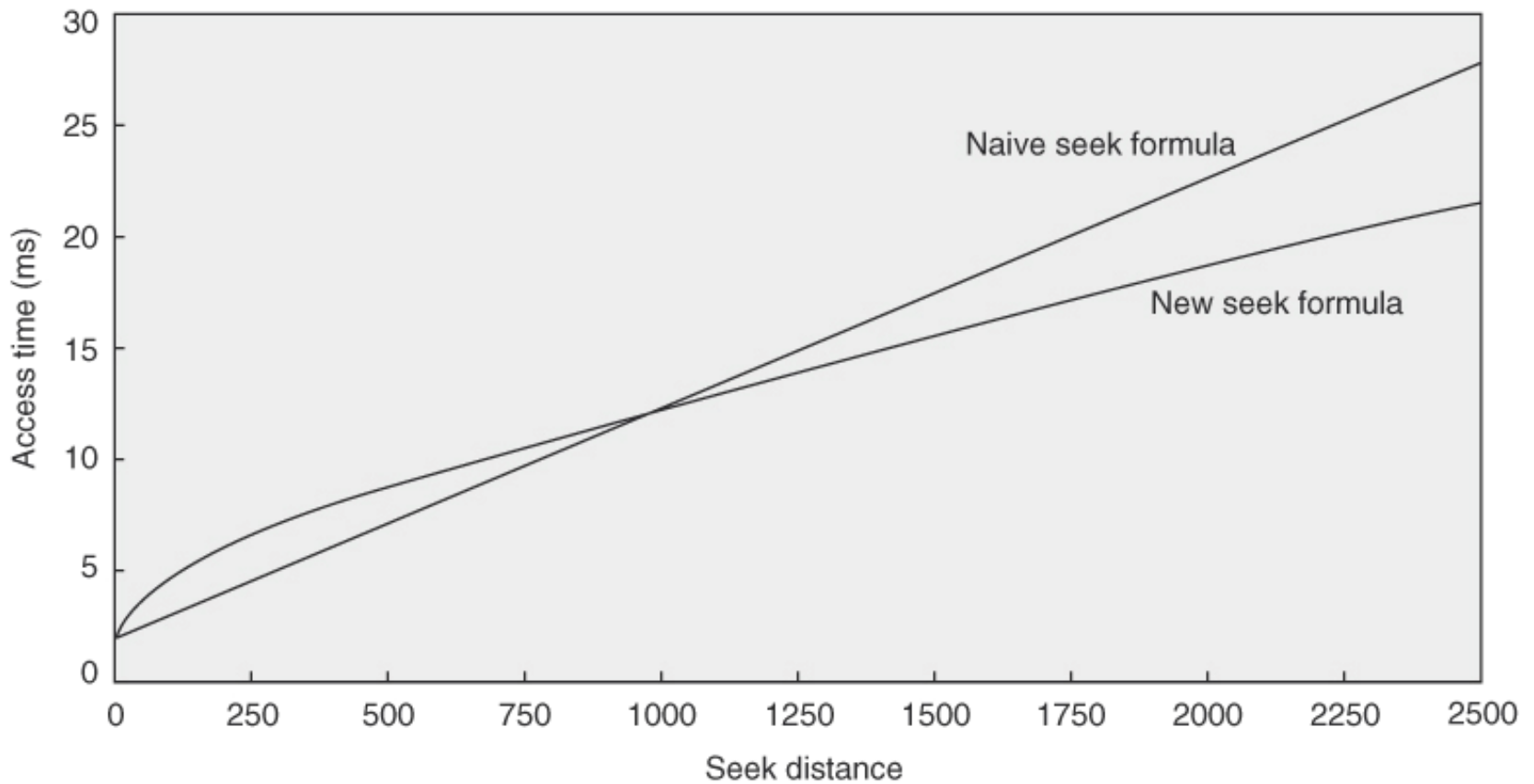
**Figure D.22 Example showing OS versus disk schedule accesses, labeled host-ordered versus drive-ordered.** The former takes 3 revolutions to complete the 4 reads, while the latter completes them in just 3/4 of a revolution. (From Anderson [2003].)

The figure shows a graph with Access time (ms) on the y-axis (0 to 30) and Seek distance on the x-axis (0 to 2500). Two curves are labeled "Naive seek formula" and "New seek formula."

Below the graph:

$$a = \frac{-10 \times \text{Time}_{min} + 15 \times \text{Time}_{avg} - 5 \times \text{Time}_{max}}{3 \times \sqrt{\text{Number of cylinders}}} \qquad b = \frac{7 \times \text{Time}_{min} - 15 \times \text{Time}_{avg} + 8 \times \text{Time}_{max}}{3 \times \text{Number of cylinders}} \qquad c = \text{Time}_{min}$$

**Figure D.23 Seek time versus seek distance for sophisticated model versus naive model.** Chen and Lee [1995] found that the equations shown above for parameters *a*, *b*, and *c* worked well for several disks.
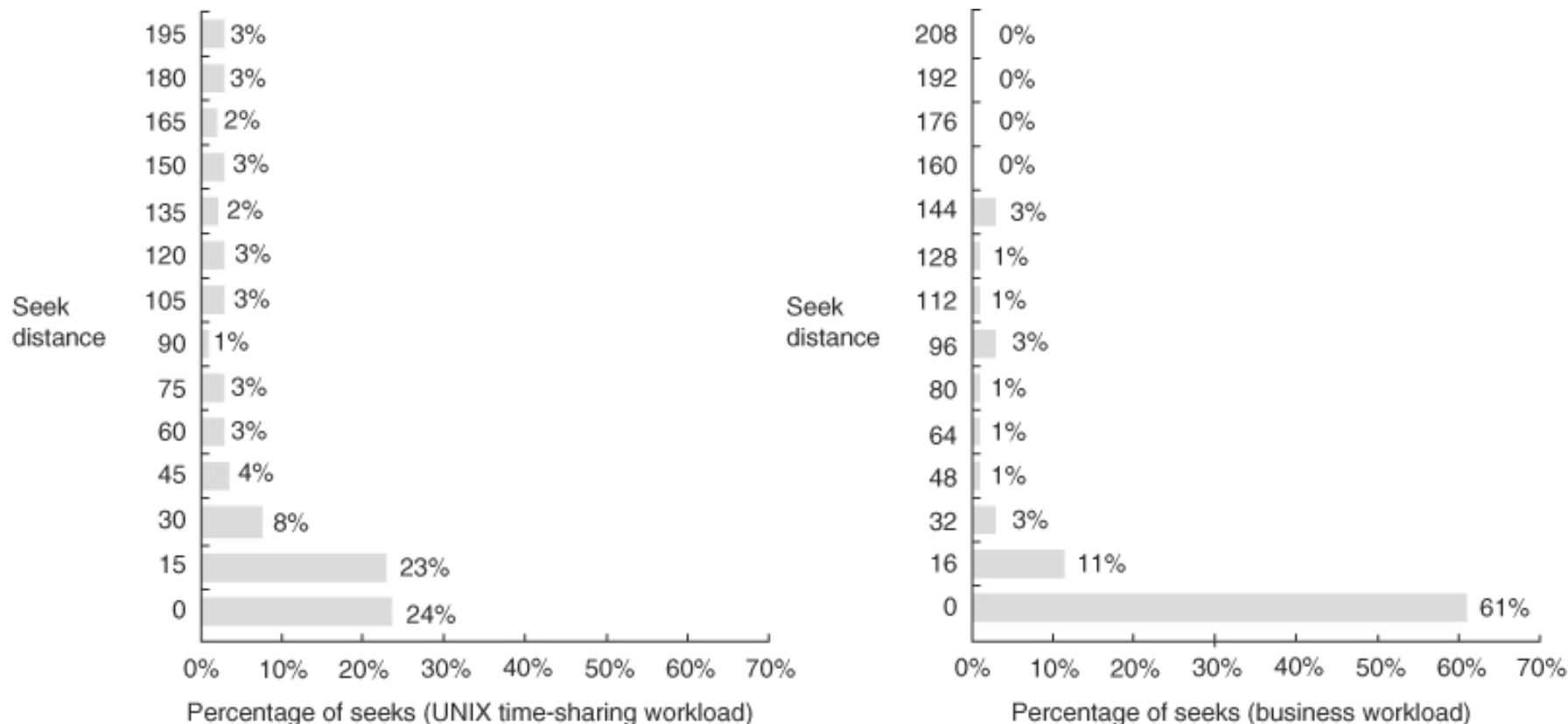
**Figure D.24 Sample measurements of seek distances for two systems.** The measurements on the left were taken on a UNIX time-sharing system. The measurements on the right were taken from a business-processing application in which the disk seek activity was scheduled to improve throughput. Seek distance of 0 means the access was made to the same cylinder. The rest of the numbers show the collective percentage for distances between numbers on the *y*-axis. For example, 11% for the bar labeled 16 in the business graph means that the percentage of seeks between 1 and 16 cylinders was 11%. The UNIX measurements stopped at 200 of the 1000 cylinders, but this captured 85% of the accesses. The business measurements tracked all 816 cylinders of the disks. The only seek distances with 1% or greater of the seeks that are not in the graph are 224 with 4%, and 304, 336, 512, and 624, each having 1%. This total is 94%, with the difference being small but nonzero distances in other categories. Measurements courtesy of Dave Anderson of Seagate.
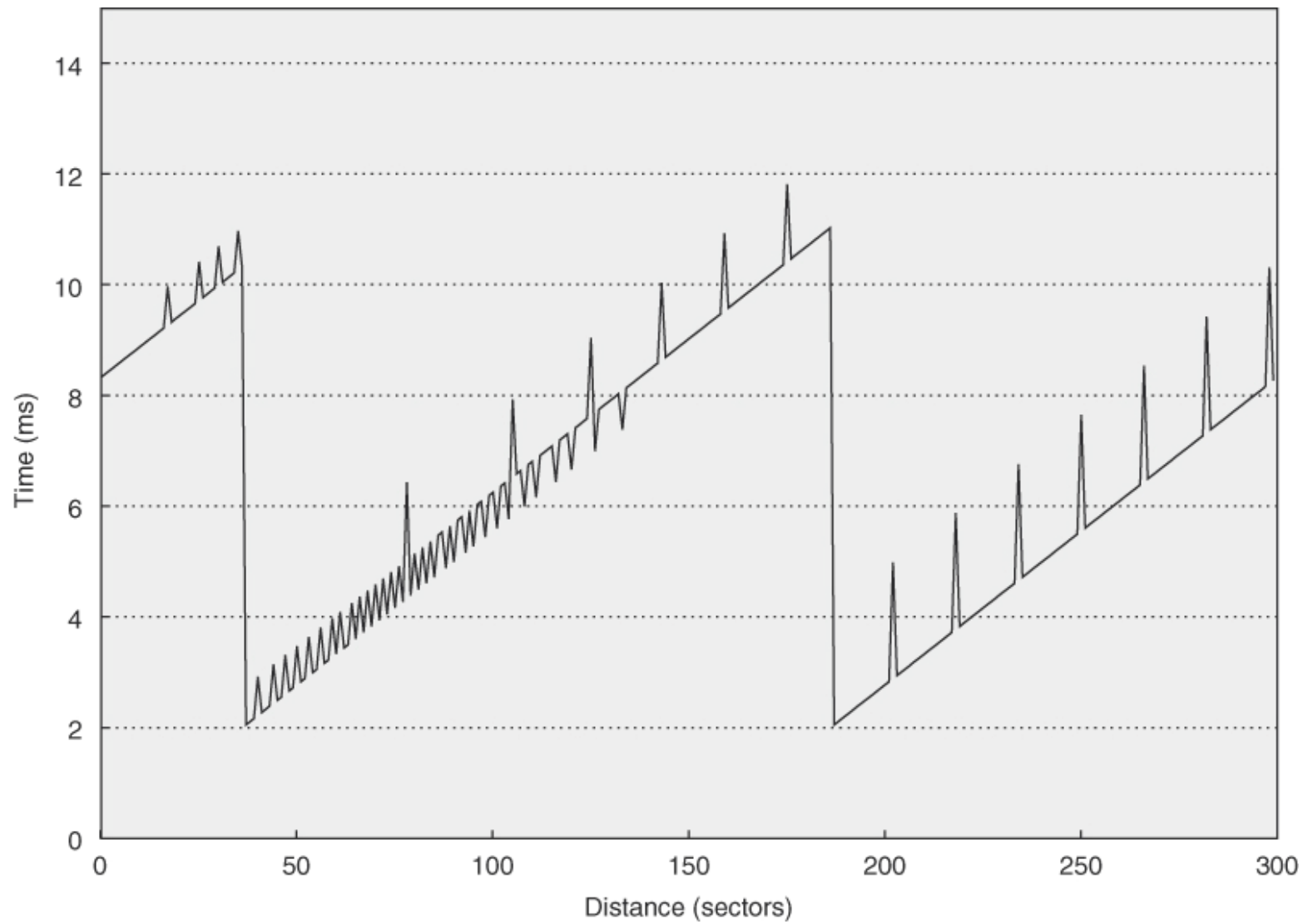
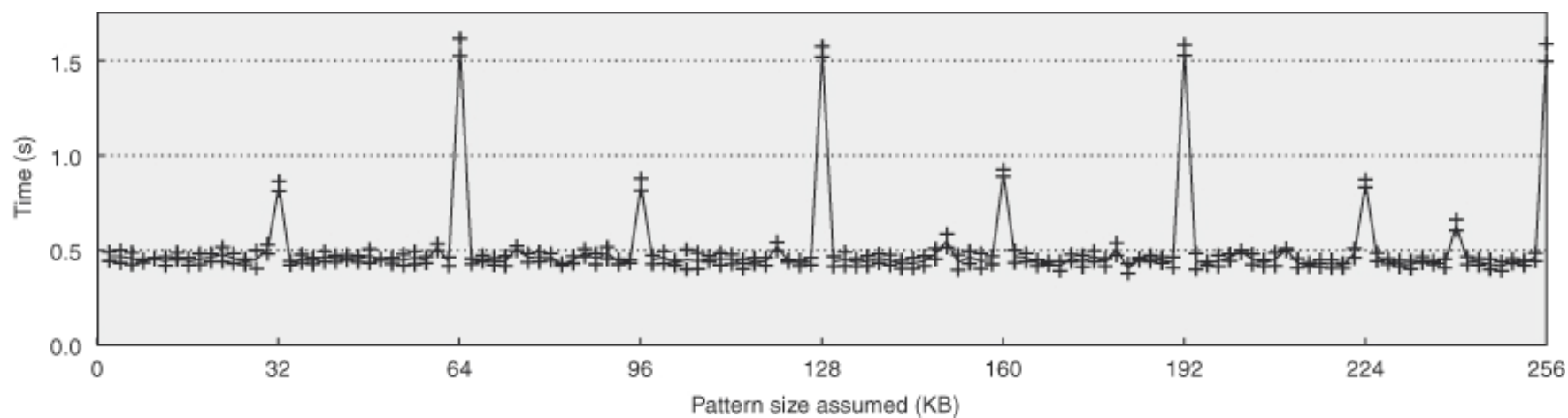**Figure D.25 Results from running Skippy on Disk Alpha.**

18

**Figure D.26 Results from running the pattern size algorithm of Shear on a mock storage system.**
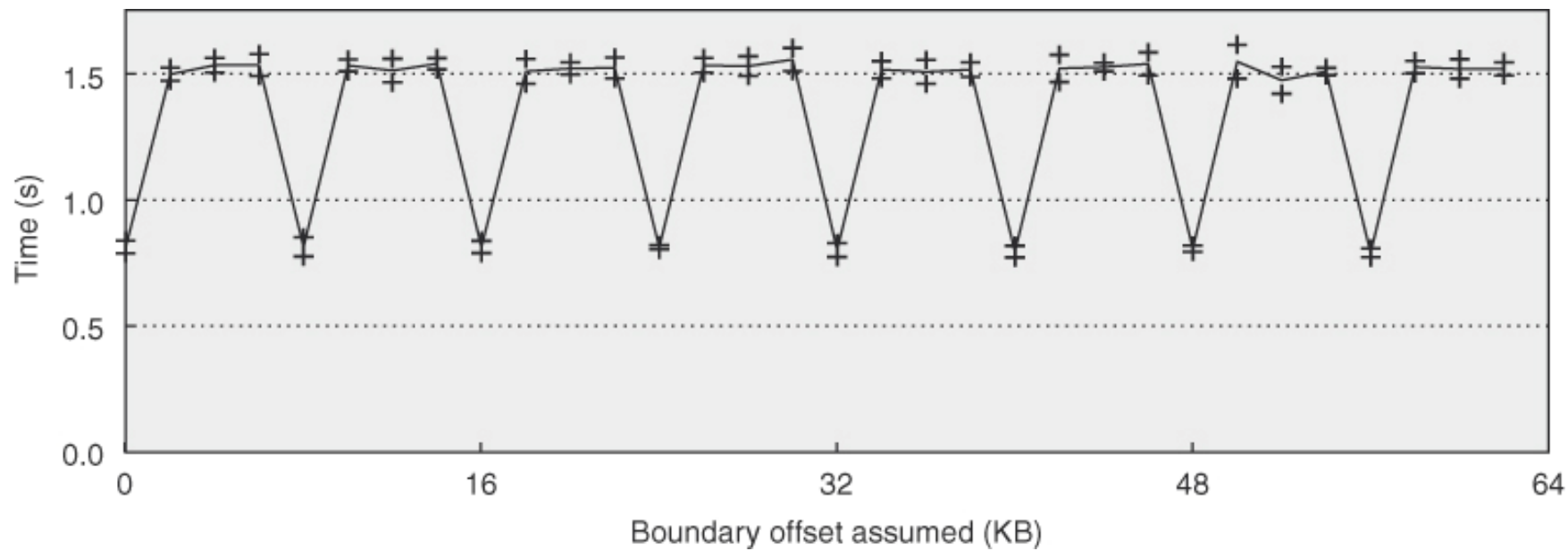
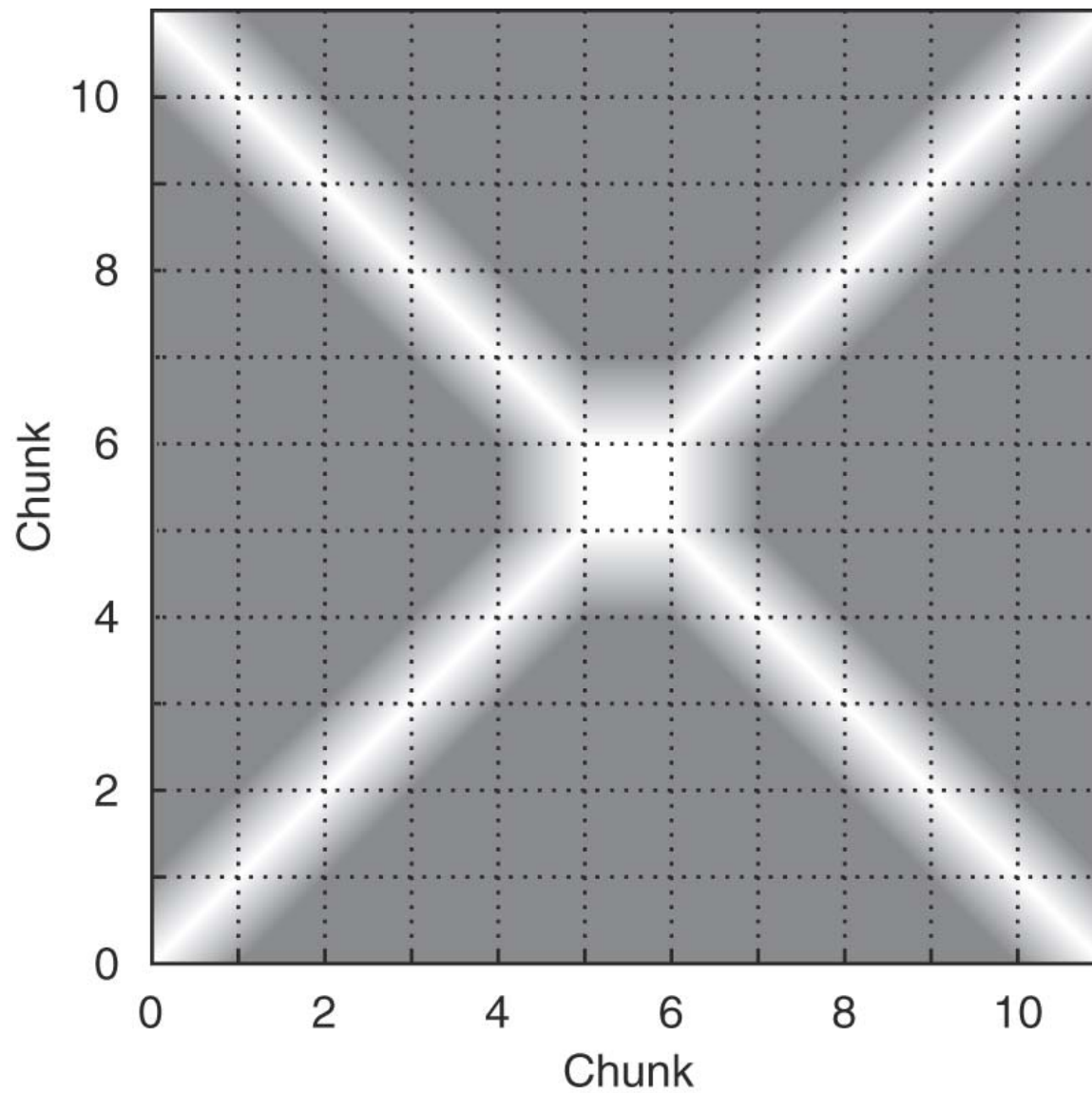**Figure D.27 Results from running the chunk size algorithm of Shear on a mock storage system.**

**Figure D.28 Results from running the layout algorithm of Shear on a mock storage system.**

Parity: RAID 5 Left–Asymmetric, stripe = 16, pattern = 48

**Figure D.29 A storage system with four disks, a chunk size of four 4 KB blocks, and using a RAID 5 Left-Asymmetric layout.** Two repetitions of the pattern are shown.