

1 Design and Evaluation of a New Approach to RAID-0 Scaling

2 GUANGYAN ZHANG and WEIMIN ZHENG, Tsinghua University

3 KEQIN LI, State University of New York and Tsinghua University

4 Scaling up a RAID-0 volume with added disks can increase its storage capacity and I/O bandwidth simul-
 5 taneously. For preserving a round-robin data distribution, existing scaling approaches require all the data
 6 to be migrated. Such large data migration results in a long redistribution time as well as a negative impact
 7 on application performance. In this article, we present a new approach to RAID-0 scaling called FastScale.
 8 First, FastScale minimizes data migration, while maintaining a uniform data distribution. It moves only
 9 enough data blocks from old disks to fill an appropriate fraction of new disks. Second, FastScale optimizes
 10 data migration with access aggregation and lazy checkpoint. Access aggregation enables data migration
 11 to have a larger throughput due to a decrement of disk seeks. Lazy checkpoint minimizes the number of
 12 metadata writes without compromising data consistency. Using several real system disk traces, we evalu-
 13 ate the performance of FastScale through comparison with SLAS, one of the most efficient existing scaling
 14 approaches. The experiments show that FastScale can reduce redistribution time by up to 86.06% with
 15 smaller application I/O latencies. The experiments also illustrate that the performance of RAID-0 scaled
 16 using FastScale is almost identical to, or even better than, that of the round-robin RAID-0.

17 Categories and Subject Descriptors: D.4.2 [Operating Systems]: Storage Management; H.3.2 [Informa-
 18 tion Storage and Retrieval]: Information Storage

19 General Terms: Algorithms, Design, Experimentation, Management, Performance

20 Additional Key Words and Phrases: Access aggregation, data migration, lazy checkpoint, RAID scaling

21 ACM Reference Format:

22 Zhang, G., Zheng, W., and Li, K. 2013. Design and evaluation of a new approach to RAID-0 scaling. *ACM*
 23 *Trans. Storage* 9, 4, Article 11 (November 2013), 31 pages.

24 DOI: <http://dx.doi.org/10.1145/2491054>

25 1. INTRODUCTION

26 1.1. Motivation

27 Redundant Array of Inexpensive Disks (RAID) was proposed to achieve high perfor-
 28 mance, large capacity, and data reliability, while allowing a RAID volume to be man-
 29 aged as a single device [Patterson et al. 1988]. This goal is achieved via disk striping
 30 and rotated parity. RAID has been well studied and widely used in high bandwidth
 31 and space-demanding areas. As user data increase and computing power is enhanced,
 32 applications often require larger storage capacity and higher I/O throughput. The scal-
 33 ability issue of RAID has become a main concern. To provide the required capacity

This work was supported by the National Natural Science Foundation of China under Grants 60903183, 61170008, and 61272055, the National High Technology Research and Development 863 Program of China under Grant 2013AA01A210, and the National Grand Fundamental Research 973 Program of China under Grant No. 2014CB340402.

Authors' addresses: G. Zhang and W. Zheng, Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China; email: {gyzh, zwm-dcs}@tsinghua.edu.cn; K. Li, Department of Computer Science, State University of New York, New Paltz, New York 12561, USA; email: lik@newpaltz.edu.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2013 ACM 1553-3077/2013/11-ART11 \$15.00

DOI: <http://dx.doi.org/10.1145/2491054>

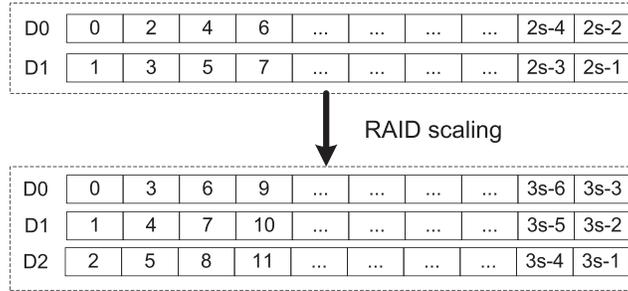


Fig. 1. RAID scaling from 2 disks to 3 using traditional approaches. All data blocks except blocks 0 and 1 have to be migrated.

and/or bandwidth, one solution is to add new disks to a RAID volume. Such disk addition is termed *RAID scaling*.

To regain uniform data distribution in all disks including the old and the new ones, RAID scaling requires certain blocks to be moved onto added disks. Furthermore, in today's server environments, many applications (e.g., e-business, scientific computation, and Web servers) access data constantly. The cost of downtime is extremely high [Patterson 2002], giving rise to the necessity of online and real-time scaling.

Traditional approaches to RAID scaling [Brown 2006; Gonzalez and Cortes 2004; Zhang et al. 2007] are restricted by preserving the round-robin order after adding disks. Let \mathbb{N} be the set of nonnegative integers. The addressing algorithm $f_i(x) : \mathbb{N} \rightarrow \mathbb{N} \times \mathbb{N}$ can be expressed as follows for the i th scaling operation:

$$f_i(x) = (d, b) = (x \bmod N_i, x/N_i), \quad (1)$$

where block $b = x/N_i$ of disk $d = x \bmod N_i$, is the location of logical block x , and N_i gives the total number of disks. As far as RAID scaling from N_{i-1} disks to N_i is concerned, since $N_i \neq N_{i-1}$, we have $f_i(x) \neq f_{i-1}(x)$ for almost every block x . Supposing $f_i(x) = f_{i-1}(x)$, we have $x = b \times N_{i-1} + d$ and $x = b \times N_i + d$, which implies that $b \times (N_i - N_{i-1}) = 0$. Since $N_i \neq N_{i-1}$, we have $b = 0$. In other words, only the data blocks in the first stripe ($b = 0$) are not moved. As an example, Figure 1 illustrates the data distributions before and after RAID scaling from 2 disks to 3. We can see that all data blocks except blocks 0 and 1 are moved during this scaling.

Suppose each disk consists of s data blocks. Let r_1 be the fraction of data blocks to be migrated. We have

$$r_1 = \frac{N_{i-1} \times s - N_{i-1}}{N_{i-1} \times s} = 1 - \frac{1}{s}. \quad (2)$$

Since s is very large, we have $r_1 \approx 100\%$. This indicates that almost 100 percent of data blocks have to be migrated no matter what the numbers of old disks and new disks are. There have been some efforts concentrating on optimization of data migration [Gonzalez and Cortes 2004; Zhang et al. 2007]. They improve the performance of RAID scaling to a certain extent, but do not completely overcome the limitation of large data migration.

The most intuitive method to reduce data migration is the semi-RR algorithm [Goel et al. 2002]. It requires a block to be moved only if the resulting disk number is one of the new disks. The algorithm can be expressed as follows for the i th scaling operation:

$$g_i(x) = \begin{cases} g_{i-1}(x), & \text{if } (x \bmod N_i) < N_{i-1}; \\ f_i(x), & \text{otherwise.} \end{cases} \quad (3)$$

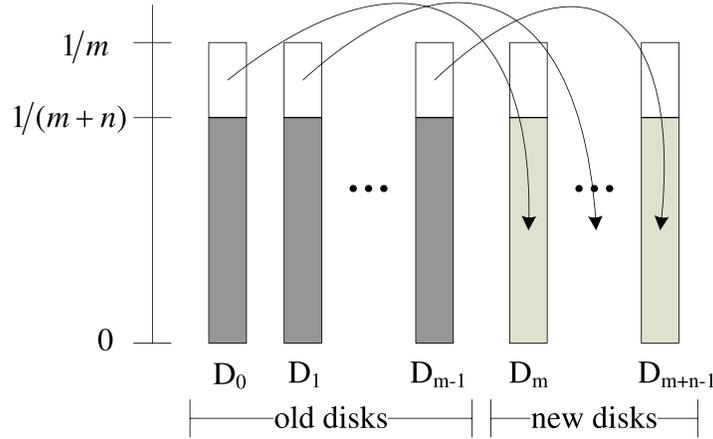


Fig. 2. Data migration using FastScale. A minimum amount of data blocks are moved from old disks to new disks to regain a uniform distribution, while no datum is migrated among old disks.

67 Semi-RR reduces data migration significantly. Unfortunately, it does not guarantee
 68 uniform distribution of data blocks after successive scaling operations (see Section 2.6).
 69 This will deteriorate the initial evenly distributed load.

70 It is clear that it is still an open problem as to whether there exists a RAID scaling
 71 method that is able to maintain a uniform and balanced load distribution by perform-
 72 ing the minimum amount of data migration.

73 1.2. Our Contributions

74 In this article, we propose a novel approach called *FastScale* to redistribute data for
 75 RAID-0 scaling. We have made three significant contributions in developing FastScale.

76 The first contribution is that FastScale accelerates RAID-0 scaling by minimizing
 77 data migration. As shown in Figure 2, FastScale moves data blocks from old disks to
 78 new disks just enough for preserving the uniformity of data distribution, while not
 79 migrating data among old disks. Before performing RAID scaling from m disks to
 80 $m + n$, the old disks hold $m \times s$ data blocks and the new disks hold no data. After
 81 RAID scaling, each disk, either old or new, holds $1/(m + n)$ of the total data to meet
 82 the uniform distribution requirement. Without loss of generality, let r_2 be the fraction
 83 of data blocks to be migrated. We have

$$84 \quad r_2 = \frac{(N_{i-1} \times s) \times \frac{1}{N_i} \times (N_i - N_{i-1})}{N_{i-1} \times s} = \frac{N_i - N_{i-1}}{N_i}. \quad (4)$$

85 For instance, for RAID scaling from 3 disks to 4, we have $r_2 = (4 - 3)/4 = 25\%$. To
 86 regain a uniform distribution, data migration from old disks to new ones is necessary.
 87 Consequently, the migration fraction r_2 of FastScale reaches the lower bound r^* of the
 88 migration fraction, where $r^* = 1 - N_{i-1}/N_i$. In other words, FastScale succeeds in
 89 minimizing data migration for RAID-0 scaling.

90 We design an elastic addressing function through which the location of one block can
 91 be easily computed without any lookup operation. By using this function, FastScale
 92 changes only a fraction of the data layout while preserving the uniformity of data
 93 distribution. FastScale has several unique features, which are listed as follows.

- 94 — FastScale maintains a uniform data distribution after each RAID scaling.
- 95 — FastScale minimizes the amount of data to be migrated during each RAID scaling.

- 96 — FastScale preserves simple management of data due to deterministic placement.
- 97 — FastScale can sustain the three preceding features after multiple disk additions.

98 The second contribution is that FastScale exploits special physical properties to opti-
 99 mize online data migration. First, it uses aggregated accesses to improve the efficiency
 100 of data migration. Second, it records data migration lazily to minimize the number of
 101 metadata updates while ensuring data consistency.

102 The third contribution is that FastScale has significant performance improvement.
 103 We implement a detailed simulator that uses DiskSim as a worker module to simulate
 104 disk accesses. Under several real system workloads, we evaluate the performance of
 105 traditional approaches and the FastScale approach. The experimental results demon-
 106 strate the following results.

- 107 — Compared with SLAS, one of the most efficient traditional approaches, FastScale
 108 shortens redistribution time by up to 86.06% with smaller maximum response time
 109 of user I/Os.
- 110 — The performance of RAID scaled using FastScale is almost identical to, or even bet-
 111 ter than, that of the round-robin RAID.

112 In this article, we only describe our solution for RAID-0, i.e., striping without par-
 113 ity. The solution can also work for RAID-10 and RAID-01. Therefore, FastScale can
 114 be used in disk arrays, logical volume managers, and file systems. Although we do
 115 not handle RAID-4 and RAID-5, we believe that our method provides a good starting
 116 point for efficient scaling of RAID-4 and RAID-5 arrays. We will report this part of our
 117 investigation in a future paper.

118 1.3. Differences from Our Prior Work

119 This article is based on our prior work presented at the 9th USENIX Conference on
 120 File and Storage Technologies (FAST'11) [Zheng and Zhang 2011]. In this article, the
 121 following new and important materials beyond the earlier version are provided.

- 122 — We improve the addressing algorithm of FastScale, especially change the way that
 123 a block newly added after the last scaling is placed. The addressing algorithm de-
 124 scribes how FastScale maps a logical address of a RAID-0 array to a physical address
 125 of a member disk. Furthermore, the presentation of the two examples is also revised
 126 accordingly.
- 127 — We design the demapping algorithm of FastScale. In many cases, it is also required
 128 to map a physical address to a logical address. The demapping algorithm can be
 129 used to provide such a mapping. The goal of improving the addressing algorithm is
 130 to enable designing the demapping algorithm and to make it simple.
- 131 — We formally prove that FastScale satisfies all three requirements for an ideal ap-
 132 proach to RAID-0 scaling.
- 133 — We also perform more experiments in different cases to make the performance eval-
 134 uation of FastScale more adequate and more convincing.
- 135 — Finally, we add some new materials to make the motivation of our work clearer and
 136 to help the reader better understand how FastScale works.

137 1.4. Article Organization

138 The rest of the article is organized as follows. In Section 2, we formally define the prob-
 139 lem to be addressed in the article, give illustrative and motivating examples, develop
 140 our addressing and demapping algorithms, and prove the properties of FastScale. In
 141 Section 3, we present the optimization techniques used in FastScale, i.e., access aggre-
 142 gation and lazy checkpoint. In Section 4, we demonstrate our experimental results to
 143 compare the performance of FastScale with that of the existing solutions. In Section 5,

144 we review related work in scaling deterministic and randomized RAID. We conclude
145 the article in Section 6.

146 2. MINIMIZING DATA MIGRATION

147 2.1. Problem Statement

148 For disk addition into a RAID-0 array, it is desirable to ensure an even load distribution
149 on all the disks and minimal block movement. Since the location of a block may be
150 changed during a scaling operation, another objective is to quickly compute the current
151 location of a block.

152 As far as the i th RAID scaling operation from N_{i-1} disks to N_i is concerned, sup-
153 pose that each disk consists of s data blocks. Before this scaling operation, there are
154 $N_{i-1} \times s$ blocks stored on N_{i-1} disks. To achieve these objectives, the following three
155 requirements should be satisfied for RAID scaling.

- 156 — *Requirement 1 (Uniform Data Distribution)*. After this scaling operation, the ex-
157 pected number of data blocks on each one of the N_i disks is $(N_{i-1} \times s)/N_i$, so as to
158 maintain an even load distribution.
- 159 — *Requirement 2 (Minimal Data Migration)*. During this scaling operation, the ex-
160 pected number of data blocks to be moved is $(N_{i-1} \times s) \times (N_i - N_{i-1})/N_i$.
- 161 — *Requirement 3 (Fast Data Addressing)*. After this scaling operation, the location of
162 a block is computed by an algorithm with low space and time complexities for all
163 original and unmoved data, original and migrated data, and new data.

164 2.2. Examples of RAID Scaling Using FastScale

165 *Example 1.* To understand how the FastScale algorithm works and how it satisfies
166 all three requirements, we take RAID scaling from 3 disks to 5 as an example. As
167 shown in Figure 3, one RAID scaling process can be divided into two stages logically,
168 i.e., data migration and data filling. In the first stage, a fraction of existing data blocks
169 are migrated to new disks. In the second stage, new data are filled into the RAID con-
170 tinuously. Actually, the two stages, data migration and data filling, can be overlapped
171 in time.

172 For the RAID scaling, each 5 consecutive locations in one disk are grouped into one
173 *segment*. For the 5 disks, 5 segments with the same physical address are grouped into
174 one *region*. Locations on all disks with the same block number form a *column* or a
175 *stripe*. In Figure 3, different regions are separated by wavy lines. For different regions,
176 the ways for data migration and data filling are completely identical. Therefore, we
177 will focus on one region, and let $s = 5$ be the number of data blocks in one segment.

178 In a region, all data blocks within a parallelogram will be moved. The base of the
179 parallelogram is 2, and the height is 3. In other words, 2 data blocks are selected from
180 each old disk and migrated to new disks. The 2 blocks are sequential—the start ad-
181 dress is the disk number *disk.no*. Figure 3 depicts the moving trace of each migrating
182 block. For one moving data block, its physical disk number is changed while its phys-
183 ical block number is unchanged. As a result, the five columns of two new disks will
184 contain 1, 2, 2, 1, and 0 migrated data blocks, respectively. Here, the data block in the
185 first column will be placed on disk 3, while the data block in the fourth column will
186 be placed on disk 4. The first blocks in columns 2 and 3 are placed on disk 3, and the
187 second blocks in columns 2 and 3 are placed on disk 4. Thus, each new disk has 3 data
188 blocks.

189 After data migration, each disk, either old or new, has 3 data blocks. That is to say,
190 FastScale regains a uniform data distribution. The total number of data blocks to be
191 moved is $2 \times 3 = 6$. This reaches the minimal number of moved blocks in each region,

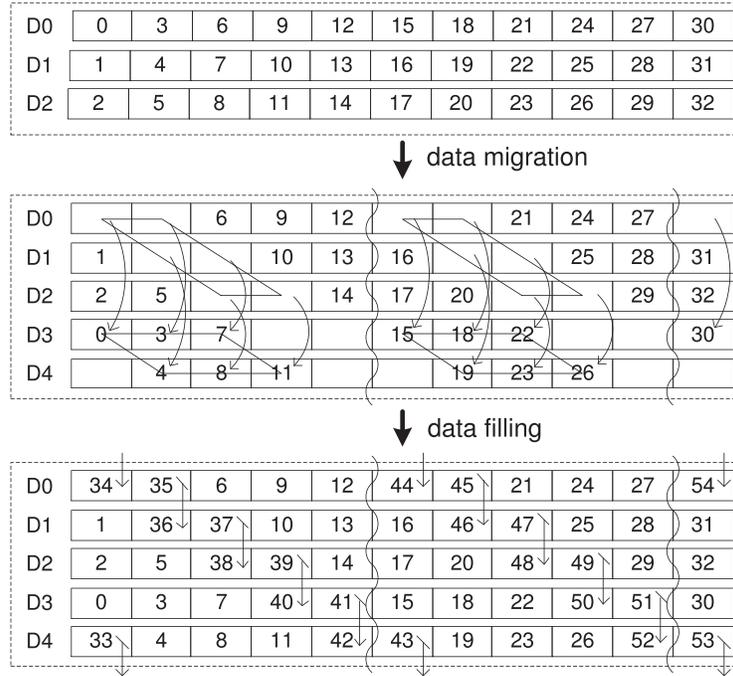


Fig. 3. RAID scaling from 3 disks to 5 using FastScale, where $m \geq n$.

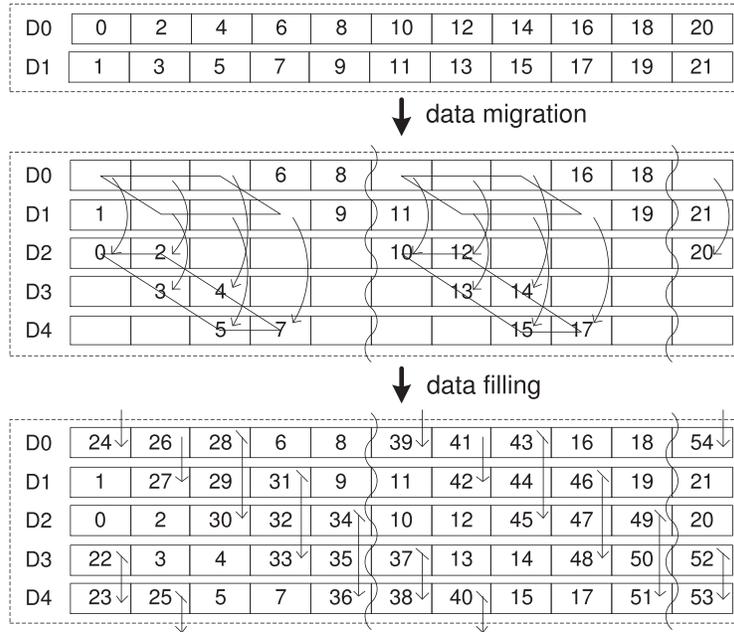
192 i.e., $(N_{i-1} \times s) \times (N_i - N_{i-1})/N_i = (3 \times 5) \times (5 - 3)/5 = 6$. We can claim that the RAID
 193 scaling using FastScale can satisfy Requirement 1 and Requirement 2.

194 Let us examine whether FastScale can satisfy Requirement 3, i.e., fast data address-
 195 ing. To consider how one logical data block is addressed, we divide the data space in
 196 the RAID into three categories: original and unmoved data, original and migrated
 197 data, and new data. The conclusion can be drawn from the following description that
 198 the calculation overhead for the data addressing is very low.

- 199 — The original and unmoved data can be addressed with the original addressing
 200 method. In this example, the ordinal number of the disk that holds one block x can
 201 be calculated as $d = x \bmod 3$. Its physical block number can be calculated as $b = x/3$.
- 202 — The addressing method for original and migrated data can be obtained easily from
 203 the description about the trace of the data migration. First, we have $b = x/3$. Next,
 204 for those blocks in the first triangle, i.e., blocks 0, 3, and 4, we have $d = d_0 + 3$. For
 205 those blocks in the last triangle, i.e., blocks 7, 8, and 11, we have $d = d_0 + 2$. Here,
 206 d_0 is the original disk number of a block.
- 207 — Each region can hold $5 \times 2 = 10$ new blocks. In one region, how those new data
 208 blocks are placed is shown in Figure 3. If block x is a new block, it is the y th new
 209 block, where $y = x - 3 \times 11$. Each stripe holds 2 new blocks. So, we have $b = y/2$. For
 210 each new data block, we have $d = (((b - 1) \bmod 5) + (y \bmod 2)) \bmod 5$.

211 *Example 2.* In the preceding example, the number of old disks m and the number of
 212 new disks n satisfy the condition $m \geq n$. In the following, we examine the case when
 213 $m < n$. Take RAID scaling from 2 disks to 5 as an example, where, $m = 2$ and $n = 3$.

214 Likewise, in a region, all data blocks within a parallelogram will be moved. The
 215 base of the parallelogram is 3, and the height is 2. 3 consecutive data blocks are se-
 216 lected from each old disk and migrated to new disks. Figure 4 depicts the trace of each

Fig. 4. RAID scaling from 2 disks to 5 using FastScale, where $m < n$.

217 migrating block. Similarly, for one moving data block, only its physical disk number
 218 is changed, while its physical block number is unchanged. As a result, five columns of
 219 three new disks will have different numbers of existing data blocks: 1, 2, 2, 1, 0. Here,
 220 the data block in the first column will be placed on disk 2, while the data block in the
 221 fourth column will be placed on disk 4. Unlike the first example, the first blocks in
 222 columns 2 and 3 are placed on disks 2 and 3, respectively. Thus, each new disk has 2
 223 data blocks.

224 After data migration, each disk, either old or new, has 2 data blocks. That is to
 225 say, FastScale regains a uniform data distribution. The total number of data blocks
 226 to be moved is $3 \times 2 = 6$. This reaches the minimal number of moved blocks, $(2 \times$
 227 $5) \times (5 - 2)/5 = 6$. We can claim that the RAID scaling using FastScale can satisfy
 228 Requirement 1 and Requirement 2.

229 Let us examine whether FastScale can satisfy Requirement 3, i.e., fast data address-
 230 ing. To consider how one logical data block is addressed, we divide the data space in
 231 the RAID into three categories: original and unmoved data, original and migrated
 232 data, and new data. The conclusion can be drawn from the following description that
 233 the calculation overhead for the data addressing is very low.

- 234 — The original and unmoved data can be addressed with the original addressing
 235 method. In this example, the ordinal number of the disk holds that one block x can
 236 be calculated as $d = x \bmod 2$. Its physical block number can be calculated as $b = x/2$.
- 237 — The addressing method for original and migrated data can be easily obtained from
 238 the description about the trace of the data migration. First, we have $b = x/2$. Next,
 239 for those blocks in the first triangle, i.e., blocks 0, 2, and 3, we have $d = d_0 + 2$. For
 240 those blocks in the last triangle, i.e., blocks 4, 5, and 7, we have $d = d_0 + 3$. Here, d_0
 241 is the original disk number of a block.

242 — Each region can hold $5 \times 3 = 15$ new blocks. In one region, how those new data
 243 blocks are placed is shown in Figure 4. If block x is a new block, it is the y th new
 244 block, where $y = x - 2 \times 11$. Each stripe holds 3 new blocks. So, we have $b = y/3$. For
 245 each new data block, we have $d = (((b - 2) \bmod 5) + (y \bmod 3)) \bmod 5$.

246 Similar to the first example, we can again claim that the RAID scaling using
 247 FastScale can satisfy the three requirements.

248 2.3. The Addressing Algorithm of FastScale

249 *2.3.1. The Addressing Algorithm.* Figure 5 shows the addressing algorithm for minimiz-
 250 ing the data migration required by RAID scaling. An array N is used to record the
 251 history of RAID scaling. $N[0]$ is the initial number of disks in the RAID. After the i th
 252 scaling operation, the RAID consists of $N[i]$ disks.

253 When a RAID is constructed from scratch ($t = 0$), it is actually a round-robin RAID.
 254 The address of block x can be calculated via one division and one modular operations
 255 (lines 3–4).

256 Let us examine the t th scaling, where n disks are added into a RAID made up of m
 257 disks (lines 7–8).

- 258 (1) If block x is an original block (line 9), FastScale calculates its old address (d_0, b_0)
 259 before the t th scaling (line 10).
 - 260 — If the block (d_0, b_0) needs to be moved (line 12), FastScale changes the disk
 261 ordinal number via the `Moving()` function (line 13), while it keeps the block
 262 ordinal number unchanged (line 14).
 - 263 — If the block (d_0, b_0) does not need to be moved (line 15), FastScale keeps the
 264 disk ordinal number and the block ordinal number unchanged (lines 16–17).
- 265 (2) If block x is a new block (line 19), FastScale places it via the `Placing()` procedure
 266 (line 20).

267 *2.3.2. The Moving Function.* The code of line 12 in Figure 5 is used to decide whether a
 268 data block (d_0, b_0) will be moved during a RAID scaling. As shown in Figures 3 and
 269 4, there is a parallelogram in each region. The base of the parallelogram is n , and the
 270 height is m . If and only if the data block is within a parallelogram, it will be moved.
 271 One parallelogram mapped to disk d_0 is a line segment. Its beginning and ending
 272 columns are d_0 and $d_0 + n - 1$, respectively. If b_1 is within the line segment, block x is
 273 within the parallelogram, and therefore it will be moved.

274 Once a data block is determined to be moved, FastScale changes its disk ordinal
 275 number with the `Moving()` function given in Figure 6. As shown in Figure 7, a migrat-
 276 ing parallelogram is divided into three parts: a head triangle, a body parallelogram,
 277 and a tail triangle. How a block moves depends on which part it lies in. No matter
 278 which is bigger between m and n , the head triangle and the tail triangle keep their
 279 shapes unchanged. The head triangle will be moved by m disks (lines 3, 9), while the
 280 tail triangle will be moved by n disks (lines 5, 11). However, the body is sensitive to the
 281 relationship between m and n . The body is twisted from a parallelogram to a rectangle
 282 when $m \geq n$ (line 6), and from a rectangle to a parallelogram when $m < n$ (line 12).
 283 FastScale keeps the relative locations of all blocks in the same column.

284 *2.3.3. The Placing Procedure.* When block x is at a location newly added after the last
 285 scaling, it is addressed via the `Placing()` procedure given in Figure 8. If block x is
 286 a new block, it is the y th new block (line 1). Each stripe holds n new blocks. So we
 287 have $b = y/n$ (line 2). In stripe b , the first new block is on disk e (line 3). Block x is
 288 the r th new data block in stripe b (line 4). Therefore, the disk number of block x is

Algorithm 1: Addressing(t, N, s, x, d, b).

Input: The input parameters are t, N, s , and x , where

t : the number of scaling times;
 N : the scaling history ($N[0], N[1], \dots, N[t]$);
 s : the number of data blocks in one disk;
 x : a logical block number.

Output: The output data are d and b , where

d : the disk holding block x ;
 b : the physical block number on disk d .

```

if ( $t = 0$ ) then (1)
     $m \leftarrow N[0]$ ; //the number of initial disks (2)
     $d \leftarrow x \bmod m$ ; (3)
     $b \leftarrow x/m$ ; (4)
    return; (5)
end if; (6)
 $m \leftarrow N[t - 1]$ ; //the number of old disks (7)
 $n \leftarrow N[t] - m$ ; //the number of new disks (8)
if ( $0 \leq x \leq m \times s - 1$ ) then //an old data block (9)
    Addressing( $t - 1, N, s, x, d_0, b_0$ ); //find the address ( $d_0, b_0$ ) before the  $t$ th scaling (10)
     $b_1 \leftarrow b_0 \bmod (m + n)$ ; (11)
    if ( $d_0 \leq b_1 \leq d_0 + n - 1$ ) then //an original block to be moved (12)
         $d \leftarrow \text{Moving}(d_0, b_1, m, n)$ ; (13)
         $b \leftarrow b_0$ ; (14)
    else //an original block not to be moved (15)
         $d \leftarrow d_0$ ; (16)
         $b \leftarrow b_0$ ; (17)
    end if; (18)
else //a new data block (19)
    Placing( $x, m, n, s, d, b$ ); (20)
end if. (21)

```

Fig. 5. The Addressing algorithm used in FastScale.

289 $(e + r) \bmod (m + n)$ (line 5). The order of placing new blocks is shown in Figures 3
 290 and 4.

291 The addressing algorithm of FastScale is very simple and elegant. It requires fewer
 292 than 50 lines of C code, reducing the likelihood that a bug will cause a data block to be
 293 mapped to a wrong location.

294 2.4. The Demapping Algorithm of FastScale

295 *2.4.1. The Demapping Algorithm.* The *Addressing* algorithm describes how FastScale
 296 maps a logical address of a RAID-0 array to a physical address of a member disk. In
 297 many cases, it is also required to map a physical address to a logical address. FastScale
 298 uses the *Demapping* algorithm, shown in Figure 9, to provide a mapping from physical
 299 addresses to logical addresses. The array N records the history of RAID scaling. $N[0]$
 300 is the initial number of disks in the RAID. After the i th scaling operation, the RAID
 301 consists of $N[i]$ disks.

302 When a RAID is constructed from scratch ($t = 0$), it is actually a round-robin RAID.
 303 The logical address of the block at (d, b) can be calculated via one multiplication and
 304 one addition operation (line 3).

Algorithm 2: Moving(d_0, b_1, m, n).

Input: The input parameters are d_0, b_1, m , and n , where
 d_0 : the disk number before a block is moved;
 b_1 : the location in a region before a block is moved;
 m : the number of old disks;
 n : the number of new disks.

Output: The disk number after a block is moved.

```

if ( $m \geq n$ ) then                                     (1)
  if ( $b_1 \leq n - 1$ ) then //head                       (2)
    return  $d_0 + m$ ;                                     (3)
  if ( $b_1 \geq m - 1$ ) then //tail                       (4)
    return  $d_0 + n$ ;                                     (5)
  return  $m + n - 1 - (b_1 - d_0)$ ; //body                (6)
else                                                    (7)
  if ( $b_1 \leq m - 1$ ) then //head                       (8)
    return  $d_0 + m$ ;                                     (9)
  if ( $b_1 \geq n - 1$ ) then //tail                       (10)
    return  $d_0 + n$ ;                                     (11)
  return  $d_0 + b_1 + 1$ ; //body                          (12)
end if.                                                (13)

```

Fig. 6. The Moving function used in the Addressing algorithm.

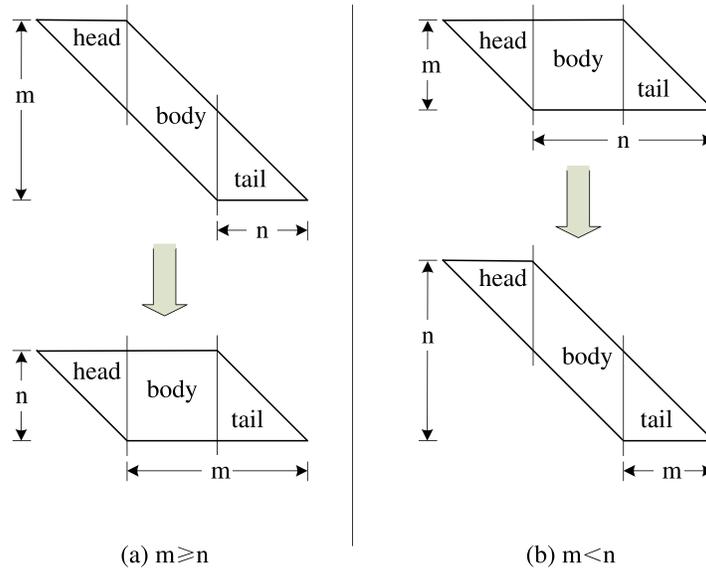


Fig. 7. The variation of data layout involved in migration.

305 Let us examine the t th scaling, where n disks are added into a RAID made up of
 306 m disks (lines 6–7). The logical address of block (d, b) can be calculated in a recursive
 307 manner.

308 — If block (d, b) is an original block and is not moved (line 9), the logical address of
 309 block (d, b) remains unchanged across the t th scaling (line 10).

Algorithm 3: $\text{Placing}(x, m, n, s, d, b)$.

Input: The input parameters are x, m, n , and s , where

- x : a logical block number;
- m : the number of old disks;
- n : the number of new disks;
- s : the number of data blocks in one disk.

Output: The output data are d and b , where

- d : the disk holding block x ;
 - b : the physical block number on disk d .
-

```

 $y \leftarrow x - m \times s;$  (1)
 $b \leftarrow y/n;$  (2)
 $e \leftarrow (b - (n - 1)) \bmod (m + n);$  (3)
 $r \leftarrow y \bmod n;$  (4)
 $d \leftarrow (e + r) \bmod (m + n).$  (5)

```

Fig. 8. The Placing procedure used in the Addressing algorithm.

Algorithm 4: $\text{Demapping}(t, N, s, d, b)$.

Input: The input parameters are t, N, s, d and b , where

- t : the number of scaling times;
- N : the scaling history ($N[0], N[1], \dots, N[t]$);
- s : the number of data blocks in one disk;
- d : the disk holding block x ;
- b : the physical block number on disk d .

Output: The output datum is data block x at location b of disk d .

```

if ( $t = 0$ ) then (1)
   $m \leftarrow N[0];$  //the number of initial disks (2)
   $x \leftarrow b \times m + d;$  (3)
  return  $x;$  (4)
end if; (5)
 $m \leftarrow N[t - 1];$  //the number of old disks (6)
 $n \leftarrow N[t] - m;$  //the number of new disks (7)
 $b_1 \leftarrow b \bmod (m + n);$  (8)
if ( $0 \leq d < m$ ) and ( $b_1 < d$  or  $b_1 > d + n - 1$ ) then //an old and unmoved data block (9)
  return  $\text{Demapping}(t - 1, N, s, d, b);$  (10)
end if; (11)
if ( $d \geq m$ ) and ( $d - m \leq b_1 \leq d - 1$ ) then //an old and moved data block (12)
   $d_0 \leftarrow \text{Demoving}(d, b_1, m, n);$  //( $d_0, b$ ) is its location before the  $t$ th scaling (13)
  return  $\text{Demapping}(t - 1, N, s, d_0, b);$  (14)
end if; (15)
return  $\text{Deplacing}(m, n, s, d, b).$  //a new data block (16)

```

Fig. 9. The Demapping algorithm used in FastScale.

310 — If block (d, b) is an original block and is moved (line 12), FastScale gets its original
 311 location (d_0, b) before the t th scaling via the $\text{Demoving}()$ function (line 13). It should
 312 be remembered that FastScale changes the disk ordinal number while keeping the
 313 block ordinal number unchanged. Then, FastScale calculates the logical address of
 314 block (d_0, b) before the t th scaling (line 14).

Algorithm 5: Demoving(d, b_1, m, n).

Input: The input parameters are d, b_1, m , and n , where
 d : the disk number after a block is moved;
 b_1 : the location in a region of the block;
 m : the number of old disks;
 n : the number of new disks.

Output: The disk number before a block is moved.

```

if ( $m \geq n$ ) then                                     (1)
  if ( $b_1 \leq n - 1$ ) then //head                       (2)
    return  $d - m$ ;                                     (3)
  if ( $b_1 \geq m - 1$ ) then //tail                     (4)
    return  $d - n$ ;                                     (5)
  return  $d + b_1 + 1 - m - n$ ; //body                  (6)
else                                                  (7)
  if ( $b_1 \leq m - 1$ ) then //head                     (8)
    return  $d - m$ ;                                     (9)
  if ( $b_1 \geq n - 1$ ) then //tail                    (10)
    return  $d - n$ ;                                     (11)
  return  $d - b_1 - 1$ ; //body                          (12)
end if.                                              (13)

```

Fig. 10. The Demoving function used in the Demapping algorithm.

315 — If block (d, b) is a new block, FastScale gets its logical address via the *Deplacing()*
 316 function (line 16).

317 The code of line 9 is used to decide whether a data block (d, b) is an old block and is
 318 not moved during this scaling. As shown in Figures 3 and 4, there is a source parallel-
 319 ogram in each region. The base of the parallelogram is n , and the height is m . A data
 320 block is moved if and only if it is within a parallelogram. One parallelogram mapped
 321 to disk d is a line segment. Its beginning and ending columns are d and $d + n - 1$,
 322 respectively. If b_1 is outside the line segment, block (d, b) is outside the parallelogram,
 323 and therefore it is not moved.

324 *2.4.2. The Demoving Function.* Likewise, the code of line 12 in Figure 9 is used to de-
 325 termine whether a data block (d, b) is an old block and has been moved during this
 326 scaling. As shown in Figures 3 and 4, there is a destination parallelogram in each re-
 327 gion. The base of the parallelogram is m , and the height is n . A data block is moved if
 328 and only if it is within a destination parallelogram. One parallelogram mapped to disk
 329 d is a line segment. Its beginning and ending columns are $d - m$ and $d - 1$, respectively.
 330 If b_1 is within the line segment, block (d, b) is within the parallelogram, and therefore
 331 it has been moved.

332 Once it is determined that a data block has been moved, FastScale gets its original
 333 disk number via the *Demoving()* function given in Figure 10. As shown in Figure 7,
 334 a migrating parallelogram is divided into three parts: a head triangle, a body paral-
 335 lelogram, and a tail triangle. How a block moves depends on which part it lies in. No
 336 matter which is bigger between m and n , the head triangle and the tail triangle keep
 337 their shapes unchanged. The head triangle will be moved by m disks (lines 3, 9), while
 338 the tail triangle will be moved by n disks (lines 5, 11). However, the body is sensitive
 339 to the relationship between m and n . The body is twisted from a parallelogram to a
 340 rectangle when $m \geq n$ (line 6), while from a rectangle to a parallelogram when $m < n$
 341 (line 12). FastScale keeps the relative locations of all blocks in the same column.

Algorithm 6: $\text{Deplacing}(m, n, s, d, b)$.

Input: The input parameters are m, n, s, d , and b , where

- m : the number of old disks;
- n : the number of new disks;
- s : the number of data blocks in one disk.
- d : the disk holding block x ;
- b : the physical block number on disk d .

Output: The output datum are the logical block number x .

```

 $e \leftarrow (b - (n - 1)) \bmod (m + n);$  (1)
 $r \leftarrow (d - e) \bmod (m + n);$  (2)
 $y \leftarrow b \times n + r;$  (3)
 $x \leftarrow m \times s + y;$  (4)
return  $x$ . (5)

```

Fig. 11. The Deplacing function used in the Demapping algorithm.

342 *2.4.3. The Deplacing Function.* If block (d, b) is at a location newly added after the t th
 343 scaling, it is addressed via the $\text{Deplacing}()$ function given in Figure 11. Each stripe
 344 holds n new blocks. In stripe b , the first new block is on disk e (line 1). Therefore, block
 345 (d, b) is the r th new block in stripe b (line 2). Since each stripe holds n new blocks,
 346 block (d, b) is the y th new data block (line 3). Furthermore, there are $m \times s$ old data
 347 blocks, and therefore, the logical address of block (d, b) is the x (line 4).

348 2.5. Properties of FastScale

349 In this section, we formally prove that FastScale satisfies all three requirements given
 350 in Section 2.1.

351 **THEOREM 2.1.** *FastScale maintains a uniform data distribution after each RAID*
 352 *scaling.*

353 **PROOF.** Assume that there are N_{i-1} old disks and $N_i - N_{i-1}$ new disks during a
 354 RAID scaling. Since each disk is divided into segments of length N_i , and a RAID vol-
 355 ume is divided into regions with the size of $N_i \times N_i$ locations, it suffices to show that
 356 FastScale maintains a uniform data distribution in each region after each RAID scal-
 357 ing. Before a RAID scaling, there are $N_i \times N_{i-1}$ blocks of data on the N_{i-1} old disks. It
 358 is clear from Figure 7 that after the scaling, each new disk holds N_{i-1} blocks of data.
 359 Since each old disk contributes $N_i - N_{i-1}$ blocks to the new disks, each old disk also
 360 holds N_{i-1} blocks of data after the scaling. Hence, the $N_i \times N_{i-1}$ data blocks in a re-
 361 gion are evenly distributed over the N_i disks, such that each disk has N_{i-1} blocks in a
 362 region. \square

363 **THEOREM 2.2.** *FastScale performs the minimum number of data migrations during*
 364 *each RAID scaling.*

365 **PROOF.** Again, it suffices to show that FastScale performs the minimum number of
 366 data migrations in each region during each RAID scaling. According to Requirement
 367 2, the minimum number of blocks to be moved is $(N_{i-1} \times s) \times (N_i - N_{i-1})/N_i$, where
 368 each old disk has s data blocks. For one region, each segment on an old disk has N_i
 369 data blocks. Therefore, the minimum number of blocks to be moved for one region is
 370 $(N_{i-1} \times N_i) \times (N_i - N_{i-1})/N_i = N_{i-1} \times (N_i - N_{i-1})$, which is exactly the number of blocks
 371 in one moving parallelogram. \square

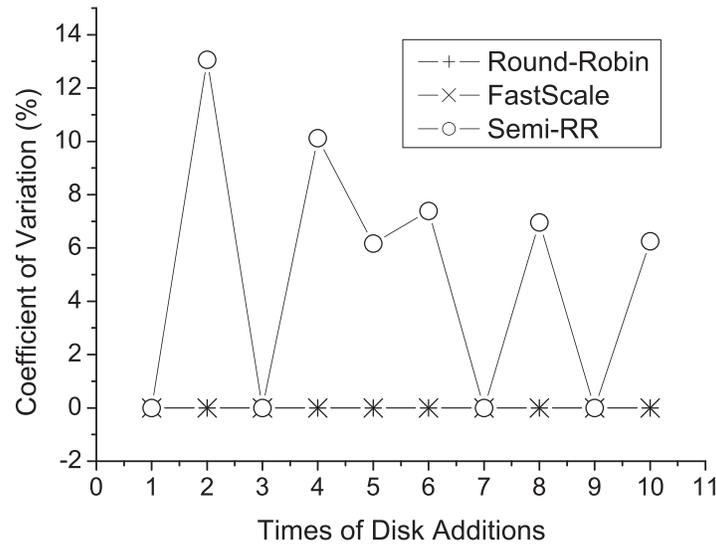


Fig. 12. Comparison of uniformity of data distributions.

372 THEOREM 2.3. *The Addressing algorithm of FastScale has time complexity $O(t)$*
 373 *after t RAID scalings.*

374 PROOF. Let $T(t)$ denote the time complexity of the Addressing algorithm. Since
 375 Addressing is a recursive algorithm, we can represent $T(t)$ by using a recurrence re-
 376 lation. First, it is clear that $T(0) = c_1$ for some constant c_1 . Next, we notice that both
 377 the Moving function and the Placing procedure take constant time. Thus, we have
 378 $T(t) \leq T(t-1) + c_2$, for all $t \geq 1$, where c_2 is some constant. Solving the recurrence
 379 relation, we get $T(t) \leq c_1 + c_2t = O(t)$. \square

380 2.6. Property Examination

381 The purpose of this experiment is to quantitatively characterize whether the FastScale
 382 algorithm satisfies the three requirements described in Section 2.1. For this purpose,
 383 we compare FastScale with the round-robin algorithm and the semi-RR algorithm.
 384 From a 4-disk array, we add one disk repeatedly for 10 times using the three algo-
 385 rithms respectively. Each disk has a capacity of 128 GB, and the size of a data block is
 386 64 KB. In other words, each disk holds 2×1024^2 blocks.

387 *2.6.1. Uniform Data Distribution.* We use the coefficient of variation of the numbers of
 388 blocks on the disks as a metric to evaluate the uniformity of data distribution across
 389 all the disks. The coefficient of variation expresses the standard deviation as a per-
 390 centage of the average. The smaller the coefficient of variation, the more uniform the
 391 data distribution. Figure 12 plots the coefficient of variation versus the number of scal-
 392 ing operations. For the round-robin and FastScale algorithms, both the coefficients of
 393 variation remain at 0 percent as the times of disk additions increase.

394 Conversely, the semi-RR algorithm causes excessive oscillation in the coefficient of
 395 variation. The maximum is even 13.06 percent. The reason for this nonuniformity is
 396 given as follows. An initial group of 4 disks causes the blocks to be placed in a round-
 397 robin fashion. When the first scaling operation adds one disk, $1/5$ of all blocks, where
 398 $(x \bmod 5) = 4$, are moved onto the new disk, i.e., disk 4. However, with another oper-
 399 ation of adding one more disk using the same approach, $1/6$ of all the blocks are not

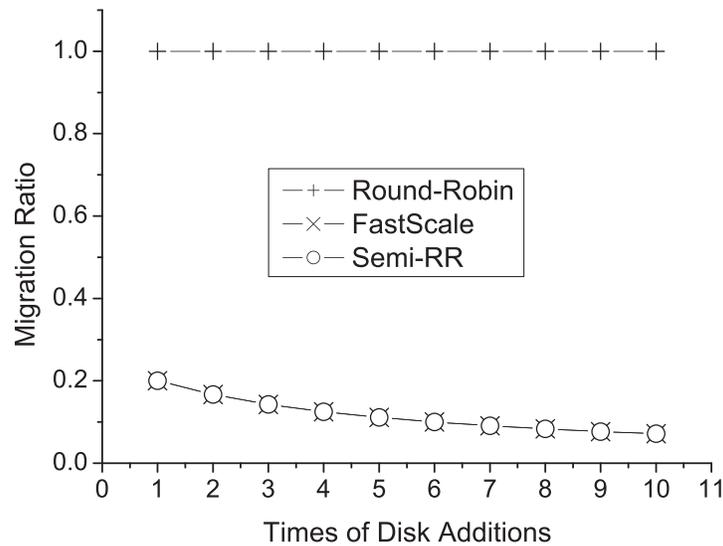


Fig. 13. Comparison of data migration ratios.

Table I. The Storage Overheads of Different Algorithms

Algorithm	Storage Overhead
Round-Robin	1
Semi-RR	t
FastScale	t

400 evenly picked from the 5 old disks and moved onto the new disk, i.e., disk 5. Only
 401 certain blocks from disks 1, 3, and 4 are moved onto disk 5 while disk 0 and disk 2
 402 are ignored. This is because disk 5 will contain blocks with logical numbers that sat-
 403 isfy $(x \bmod 6) = 5$, which are all odd numbers. The logical numbers of those blocks on
 404 disks 0 and 2, resulting from $(x \bmod 4) = 0$ and $(x \bmod 4) = 2$ respectively, are all even
 405 numbers. Therefore, blocks from disks 0 and 2 do not qualify and are not moved.

406 **2.6.2. Minimal Data Migration.** Figure 13 plots the migration fraction (the fraction of data
 407 blocks to be migrated) versus the number of scaling operations. Using the round-robin
 408 algorithm, the migration fraction is constantly 100%. This will incur a very large mi-
 409 gration cost.

410 The migration fractions using the semi-RR algorithm and using FastScale are identi-
 411 cal. They are significantly smaller than the migration fraction of using the round-robin
 412 algorithm. Another obvious phenomenon is that they decrease with the increase of the
 413 number of scaling operations. The reason for this phenomenon is described as follows.
 414 To make each new disk hold $1/N_i$ of total data, the semi-RR algorithm and FastScale
 415 move $(N_i - N_{i-1})/N_i$ of total data. N_i increases with the number of scaling operations,
 416 i.e., i . As a result, the percentage of new disks $((N_i - N_{i-1})/N_i)$ decreases. Therefore,
 417 the migration fractions using the semi-RR algorithm and FastScale decrease.

418 **2.6.3. Storage and Calculation Overheads.** When a disk array boots, it needs to obtain the
 419 RAID topology from disks. Table I shows the storage overheads of the three algorithms.
 420 The round-robin algorithm depends only on the total number of member disks. So its
 421 storage overhead is one integer. The semi-RR and FastScale algorithms depend on
 422 how many disks are added during each scaling operation. If we scale RAID t times,

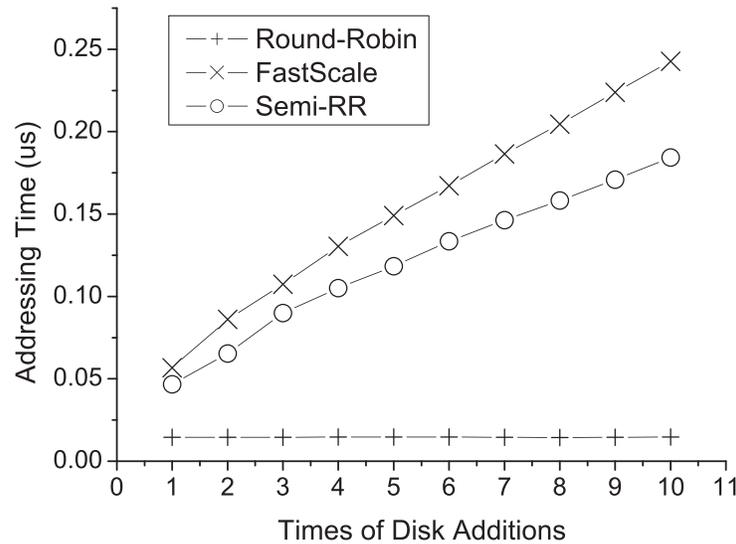


Fig. 14. Comparison of addressing times.

423 their storage overheads are t integers. Actually, the RAID scaling operation is not too
 424 frequent. It may be performed once every half a year, or even less often. Consequently,
 425 the storage overheads are very small.

426 To quantitatively characterize the calculation overheads, we run different algo-
 427 rithms to calculate the physical addresses for all data blocks on a scaled RAID. The
 428 whole addressing process is timed and then the average addressing time for each block
 429 is calculated. The testbed used in the experiment is an Intel Dual Core T9400 2.53 GHz
 430 machine with 4 GB of memory. A Windows 7 Enterprise Edition is installed. Figure 14
 431 plots the addressing time versus the number of scaling operations.

432 The round-robin algorithm has a low calculation overhead of $0.014 \mu\text{s}$ or so. The
 433 calculation overheads using the semi-RR and FastScale algorithms are close, and both
 434 take on an upward trend. Among the three algorithms, FastScale has the largest over-
 435 head. Fortunately, the largest addressing time using FastScale is $0.24 \mu\text{s}$ which is
 436 negligible compared to milliseconds of disk I/O time.

437 3. OPTIMIZING DATA MIGRATION

438 The FastScale algorithm succeeds in minimizing data migration for RAID scaling. In
 439 this section, we describe FastScale's optimizations for the process of data migration.
 440 To better understand how these optimizations work, we first give an overview of the
 441 data migration process.

442 3.1. Overview of the Scaling Process

443 Before the i th scaling operation, an addressing equation, $h_{i-1}(x)$, describes the original
 444 geometry where N_{i-1} disks serve user requests.

445 Figure 15 illustrates an overview of the migration process. FastScale uses a sliding
 446 window to describe the mapping information of a continuous segment in a RAID un-
 447 der scaling. During scaling, only data that lie within the sliding window is copied to
 448 new locations. The addressing information of the sliding window is maintained with a
 449 bitmap table, where a bit indicates whether a data block has been migrated. The size
 450 of a sliding window is exactly that of a *region*.

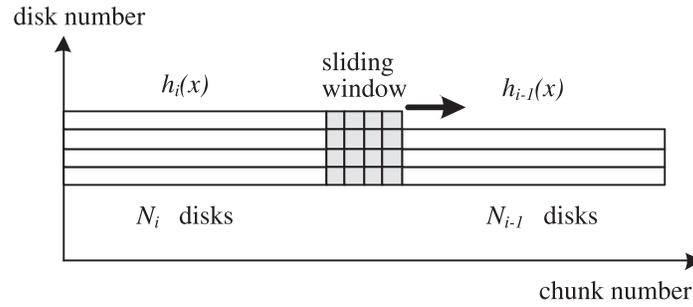


Fig. 15. An overview of the data redistribution process. As the sliding window slides ahead, the newly added disks are gradually available to serve user requests.

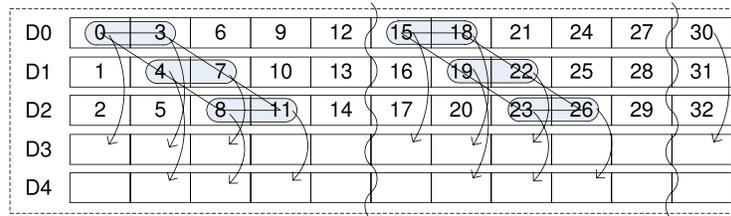


Fig. 16. Aggregate reads for RAID scaling from 3 disks to 5. Multiple successive blocks are read via a single I/O.

451 An incoming user request is mapped in one of three ways according to its logical
452 address.

- 453 — If its logical address is above the sliding window, it is mapped through the equation
- 454 $h_{i-1}(x)$, where N_{i-1} disks serve user requests.
- 455 — If its logical address is below the sliding window, it is mapped through the new
- 456 equation $h_i(x)$, where N_i disks serve user requests.
- 457 — If its logical address is in the range of the sliding window, it is mapped through the
- 458 sliding window.

459 When all the data in a sliding window are moved, the sliding window moves ahead
460 by one *region* size. In this way, the newly added disks are gradually available to serve
461 user requests. When data redistribution is completed, the new addressing equation,
462 $h_i(x)$, is used to describe the new geometry.

463 3.2. Access Aggregation

464 FastScale only moves data blocks from old disks to new disks, while not migrating
465 data among old disks. The data migration will not overwrite any valid data. As a re-
466 sult, data blocks may be moved in an arbitrary order. Since disk I/O performs much
467 better for large sequential accesses, FastScale accesses multiple successive blocks via a
468 single I/O.

469 Take a RAID scaling from 3 disks to 5 as an example (see Figure 16). Let us focus
470 on the first region. FastScale issues the first I/O request to read blocks 0 and 3, the
471 second request to read blocks 4 and 7, and the third request for blocks 8 and 11, simul-
472 taneously. By this means, to read all these blocks, FastScale requires only three I/Os,
473 instead of six. Furthermore, all of these 3 large-size data reads are on three disks.
474 They can be done in parallel, further increasing I/O rate.

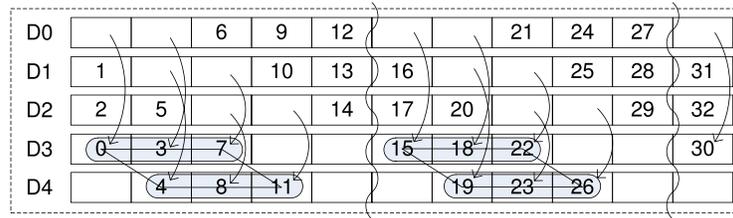


Fig. 17. Aggregate writes for RAID scaling from 3 disks to 5. Multiple successive blocks are written via a single I/O.

475 When all six blocks have been read into a memory buffer, FastScale issues the first
 476 I/O request to write blocks 0, 3, and 7, and the second I/O to write blocks 4, 8, and 11,
 477 simultaneously (see Figure 17). In this way, only two large sequential write requests
 478 are issued, as opposed to six small writes.

479 For RAID scaling from N_{i-1} disks to N_i disks, N_{i-1} reads and $N_i - N_{i-1}$ writes are
 480 required to migrate all the data in a region, i.e., $N_{i-1} \times (N_i - N_{i-1})$ data blocks.

481 Access aggregation converts sequences of small requests into fewer, larger requests.
 482 As a result, seek cost is mitigated over multiple blocks. Moreover, a typical choice of the
 483 optimal block size for RAID is 32KB or 64KB [Brown 2006; Hennessy and Patterson
 484 2003; Kim et al. 2001; Wilkes et al. 1996]. Thus, accessing multiple successive blocks
 485 via a single I/O enables FastScale to have a larger throughput. Since data densities in
 486 disks increase at a much faster rate than improvements in seek times and rotational
 487 speeds, access aggregation benefits more as technology advances.

488 3.3. Lazy Checkpoint

489 While data migration is in progress, the RAID storage serves user requests. Further-
 490 more, the coming user I/Os may be write requests to migrated data. As a result, if
 491 mapping metadata does not get updated until all blocks have been moved, data consis-
 492 tency may be destroyed. Ordered operations [Kim et al. 2001] of copying a data block
 493 and updating the mapping metadata (a.k.a., *checkpoint*) can ensure data consistency.
 494 But ordered operations cause each block movement to require one metadata write,
 495 which results in a large cost for data migration. Because metadata is usually stored
 496 at the beginning of all member disks, each metadata update causes one long seek per
 497 disk. FastScale uses lazy checkpoint to minimize the number of metadata writes with-
 498 out compromising data consistency.

499 The foundation of lazy checkpoint is described as follows. Since block copying does
 500 not overwrite any valid data, both its new replica and the original are valid after a
 501 data block is copied. In the preceding example, we suppose that blocks 0, 3, 4, 7, 8, and
 502 11 have been copied to their new locations and the mapping metadata has not been
 503 updated (see Figure 18), when the system fails. The original replicas of the six blocks
 504 will be used after the system reboots. As long as blocks 0, 3, 4, 7, 8, and 11 have not
 505 been written since being copied, the data remain consistent. Generally speaking, when
 506 the mapping information is not updated immediately after a data block is copied, an
 507 unexpected system failure only wastes some data accesses, but does not sacrifice data
 508 reliability. The only threat is the incoming of write operations to migrated data.

509 The key idea behind lazy checkpoint is that data blocks are copied to new locations
 510 continuously, while the mapping metadata is not updated onto the disks (a.k.a., *check-*
 511 *point*) until a threat to data consistency appears. We use $h_i(x)$ to describe the geometry
 512 after the i th scaling operation, where N_i disks serve user requests. Figure 19 illus-
 513 trates an overview of the migration process. Data in the moving region is copied to
 514 new locations. When a user request arrives, if its physical block address is above the

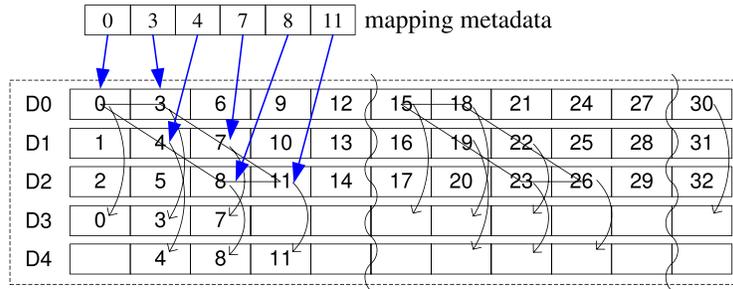


Fig. 18. If data blocks are copied to their new locations and metadata is not yet updated when the system fails, data consistency is still maintained because the data in their original locations are valid and available.

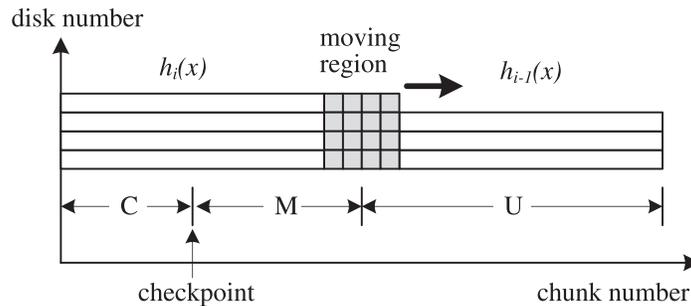


Fig. 19. Lazy updates of mapping metadata. “C”: migrated and checkpointed; “M”: migrated but not checkpointed; “U”:not migrated. Data redistribution is checkpointed only when a user write request arrives in area “M”.

515 moving region, it is mapped with $h_{i-1}(x)$; If its physical block address is below the
 516 moving region, it is mapped with $h_i(x)$. When all of the data in the current moving
 517 region are moved, the next region becomes the moving region. In this way, the newly
 518 added disks are gradually available to serve user requests. Only when a user write
 519 request arrives in the area where data have been moved and the movement has not
 520 been checkpointed, are mapping metadata updated.

521 Since one write of metadata can store multiple map changes of data blocks, lazy
 522 updates can significantly decrease the number of metadata updates, reducing the cost
 523 of data migration. Furthermore, lazy checkpoint can guarantee data consistency. Even
 524 if the system fails unexpectedly, only some data accesses are wasted. It should also be
 525 noted that the probability of a system failure is very low.

526 4. EXPERIMENTAL EVALUATION

527 The experimental results in Section 2.6 show that the semi-RR algorithm causes ex-
 528 tremely nonuniform data distribution. This will result in low I/O performance due
 529 to load imbalance. In this section, we compare FastScale with the SLAS approach
 530 [Zhang et al. 2007] through detailed experiments. SLAS, proposed in 2007, preserves
 531 the round-robin order after adding disks.

532 4.1. Simulation System

533 We use detailed simulations with several disk traces collected in real systems. The
 534 simulator is made up of a workload generator and a disk array (Figure 20). According
 535 to trace files, the workload generator initiates an I/O request at the appropriate time,
 536 so that a particular workload is induced on the disk array.

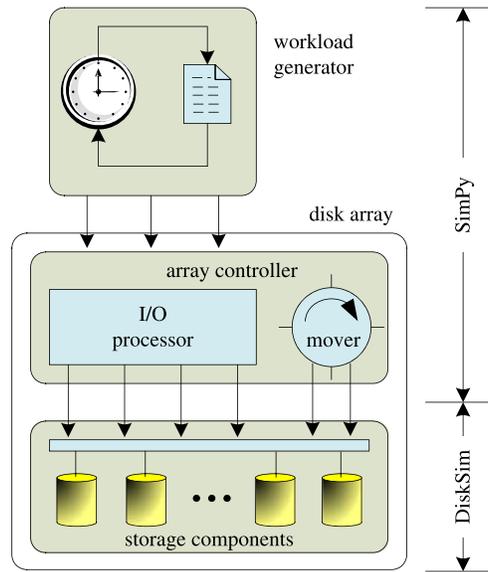


Fig. 20. A simulation system block diagram. The workload generator and the array controller were implemented in SimPy. DiskSim was used as a worker module to simulate disk accesses.

537 The disk array consists of an array controller and storage components. The array
 538 controller is logically divided into two parts: an I/O processor and a data mover. The
 539 I/O processor, according to the address mapping, forwards incoming I/O requests to
 540 the corresponding disks. The data mover reorganizes the data on the array. The mover
 541 uses on/off logic to adjust the redistribution rate. Data redistribution is throttled on
 542 detection of high application workload. Otherwise, it performs continuously. An IOPS
 543 (I/Os per second) threshold is used to determine whether an application workload
 544 is high.

545 The simulator is implemented in SimPy [Muller and Vignaux 2009] and DiskSim
 546 [Bucy et al. 2008]. SimPy is an object-oriented, process-based discrete-event simula-
 547 tion language based on standard Python. DiskSim is an efficient and accurate disk
 548 system simulator from Carnegie Mellon University and has been extensively used in
 549 various research projects studying storage subsystem architectures. The workload gen-
 550 erator and the array controller are implemented in SimPy. Storage components are
 551 implemented in DiskSim. In other words, DiskSim is used as a worker module to sim-
 552 ulate disk accesses. The simulated disk specification is that of the 15,000-RPM IBM
 553 Ultrastar 36Z15 [Hitachi 2001].

554 4.2. Workloads

555 Our experiments use the following three real system disk I/O traces with different
 556 characteristics.

- 557 — *TPC-C* traced disk accesses of the TPC-C database benchmark with 20 ware-
 558 houses [Brigham Young University, 2010.]. It was collected with one client running
 559 20 iterations.
- 560 — *Fin* is obtained from the Storage Performance Council (SPC) [UMass Trace Repos-
 561 itory 2007, Storage Performance Council 2010], a vendor-neutral standards body.
 562 The Fin trace was collected from OLTP applications running at a large financial
 563 institution. The write ratio is high.

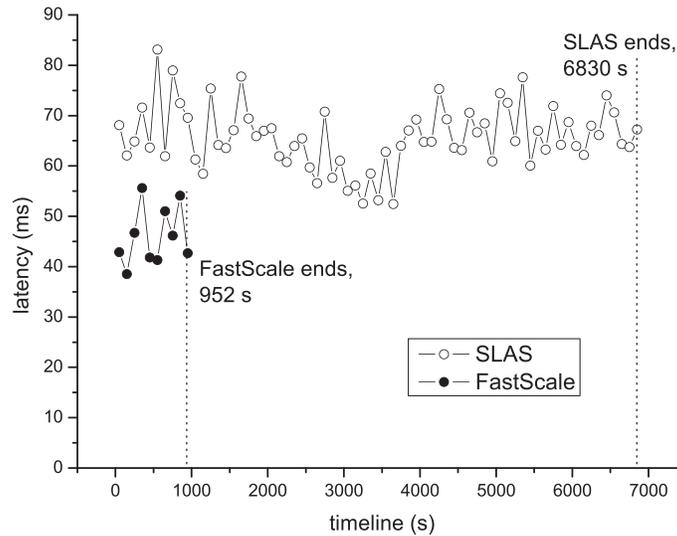


Fig. 21. Performance comparison between FastScale and SLAS under the Fin workload.

564 — *Web* is also from SPC. It was collected from a system running a Web search engine.
 565 The read-dominated Web trace exhibits strong locality in its access pattern.

566 4.3. Experiment Results

567 4.3.1. *The Scaling Efficiency.* Each experiment lasts from the beginning to the end of
 568 data redistribution for RAID scaling. We focus on comparing redistribution times and
 569 user I/O latencies when different scaling programs are running in the background.

570 In all experiments, the sliding window size for SLAS is set to 1024. Access aggrega-
 571 tion in SLAS can improve the redistribution efficiency. However, a too-large size of
 572 redistribution I/Os will compromise the I/O performance of applications. In our experi-
 573 ments, SLAS reads 8 data blocks via an I/O request.

574 The purpose of our first experiment is to quantitatively characterize the advan-
 575 tages of FastScale through a comparison with SLAS. We conduct a scaling operation
 576 of adding 2 disks to a 4-disk RAID, where each disk has a capacity of 4 GB. Each ap-
 577 proach performs with the 32KB stripe unit size under a Fin workload. The threshold
 578 of rate control is set to 100 IOPS. This parameter setup acts as the baseline for the
 579 latter experiments, from which any change will be stated explicitly.

580 We collect the latencies of all user I/Os. We divide the I/O latency sequence into
 581 multiple sections according to I/O issuing time. The time period of each section is
 582 100 seconds. Furthermore, we get a local maximum latency from each section. A lo-
 583 cal maximum latency is the maximum of I/O latency in a section. Figure 21 plots local
 584 maximum latencies using the two approaches as the time increases along the x -axis.
 585 It illustrates that FastScale demonstrates a noticeable improvement over SLAS in two
 586 metrics.

587 — First, the redistribution time using FastScale is significantly shorter than that us-
 588 ing SLAS: 952 seconds and 6830 seconds, respectively. In other words, FastScale
 589 has an 86.06% shorter redistribution time than SLAS. The main factor in
 590 FastScale’s reduction of the redistribution time is the significant decline of the
 591 amount of the data to be moved. When SLAS is used, almost 100% of the data
 592 blocks have to be migrated. However, when FastScale is used, only 33.3% of the

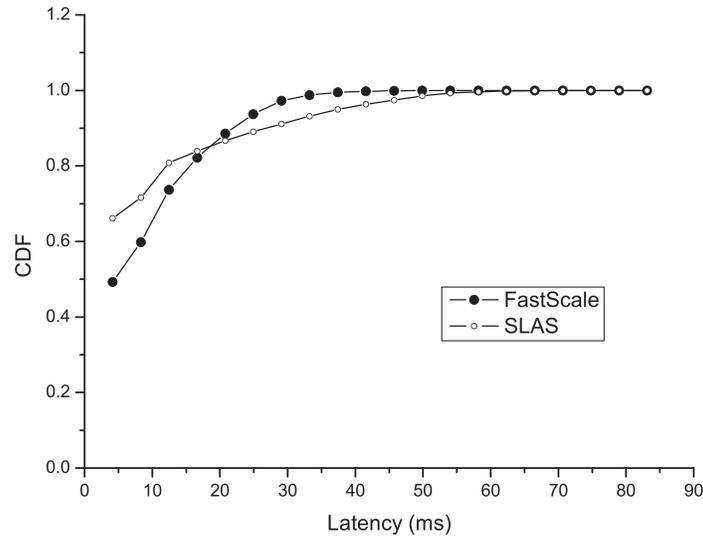


Fig. 22. Cumulative distribution of I/O latencies during data redistributions by the two approaches under the Fin workload.

593 data blocks have to be migrated. Another factor is the effective exploitation of two
 594 optimization technologies: access aggregation reduces the number of redistribution
 595 I/Os and lazy checkpoint minimizes metadata writes.
 596 — Second, local maximum latencies of SLAS are obviously longer than those of
 597 FastScale. The global maximum latency using SLAS reaches 83.12 ms while that
 598 using FastScale is 55.60 ms. This is because the redistribution I/O size using SLAS
 599 is larger than that using FastScale. For SLAS, the read size is 256 KB (8 blocks),
 600 and the write size is 192 KB (6 blocks). For FastScale, the read size is 64 KB (2
 601 blocks), and the write size is 128 KB (4 blocks). Of course, local maximum laten-
 602 cies of SLAS will be lower with a decrease in the redistribution I/O size. But the
 603 decrease in the I/O size will necessarily enlarge the redistribution time.

604 Figure 22 shows the cumulative distribution of user response times during data
 605 redistribution. To provide a fair comparison, I/Os involved in statistics for SLAS are
 606 only those issued before 952 seconds. When I/O latencies are longer than 18.65 ms,
 607 the CDF value of FastScale is greater than that of SLAS. This indicates again that
 608 FastScale has a smaller maximum response time for user I/Os than SLAS. The average
 609 latency of FastScale is close to that of SLAS: 8.01 ms and 7.53 ms respectively. It is
 610 noteworthy that due to significantly shorter data redistribution time, FastScale has a
 611 markedly smaller impact on user I/O latencies than does SLAS.

612 A factor that might affect the benefits of FastScale is the workload under which data
 613 redistribution performs. Under the TPC-C workload, we also measure the performance
 614 of FastScale and SLAS to perform the “4+2” scaling operation.

615 For the TPC-C workload, Figure 23 shows local maximum latencies versus the
 616 redistribution times for SLAS and FastScale. It once again shows the efficiency of
 617 FastScale in improving the redistribution time. The redistribution times using SLAS
 618 and FastScale are 6820 seconds and 964 seconds, respectively. That is to say, FastScale
 619 causes an improvement of 85.87% in the redistribution time. Likewise, local maximum
 620 latencies of FastScale are also obviously shorter than those of SLAS. The global maxi-
 621 mum latency using FastScale is 114.76 ms while that using SLAS reaches 147.82 ms.

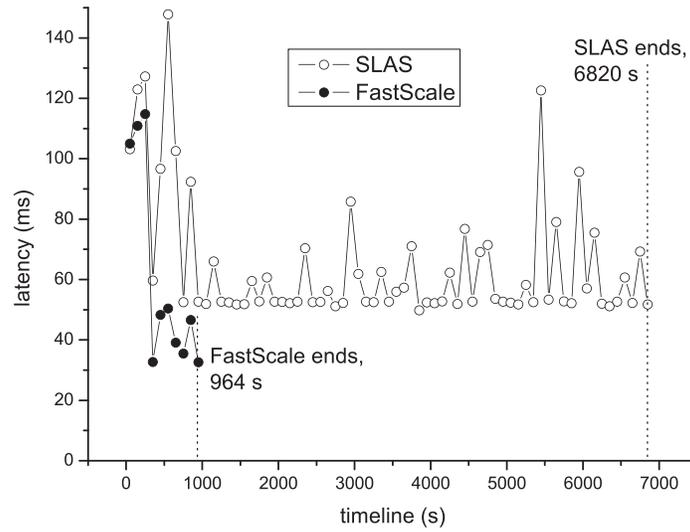


Fig. 23. Performance comparison between FastScale and SLAS under the TPC-C workload.

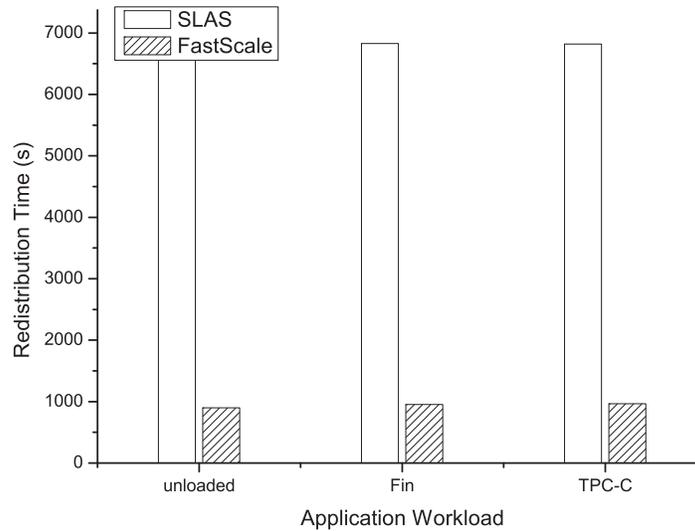


Fig. 24. Comparison of redistribution times of FastScale and SLAS under different workloads. The label “unloaded” means scaling a RAID volume offline.

622 To compare the performance of FastScale under different workloads, Figure 24
 623 shows a comparison in the redistribution time between FastScale and SLAS. For completeness,
 624 we also conducted a comparison experiment on the redistribution time with
 625 no loaded workload. To scale a RAID volume offline, SLAS uses 6802 seconds whereas
 626 FastScale consumes only 901 seconds. FastScale provides an improvement of 86.75%
 627 in the redistribution time.

628 We can draw one conclusion from Figure 24. Under various workloads, FastScale
 629 consistently outperforms SLAS by 85.87–86.75% in the redistribution time, with
 630 shorter maximum response time for user I/Os.

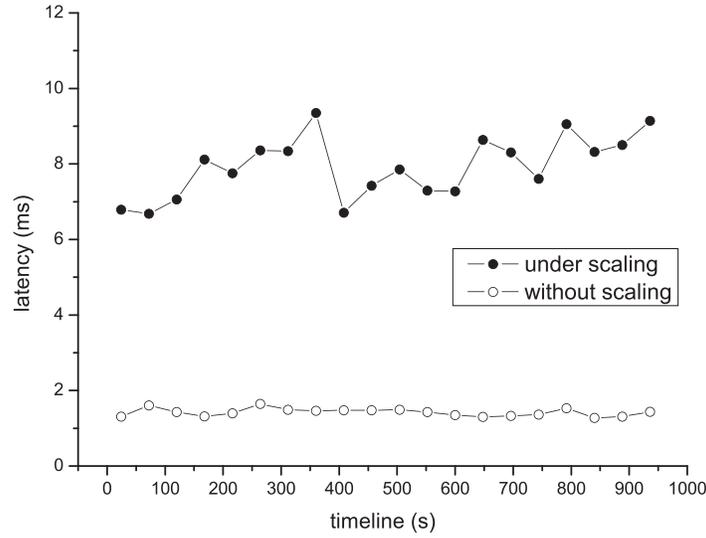


Fig. 25. Impact of RAID-0 scaling on application I/O performance under the Fin workload.

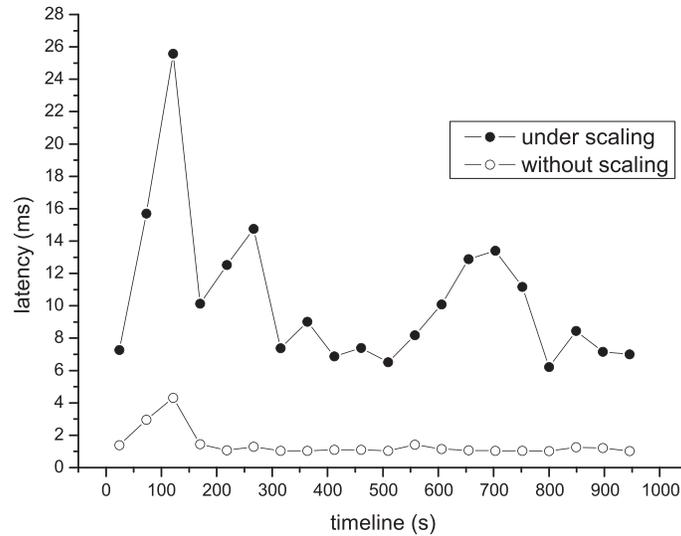


Fig. 26. Impact of RAID-0 scaling on application I/O performance under the TPC-C workload.

631 To quantitatively demonstrate how data migration affects the existing workload during
 632 RAID scaling using FastScale, we measure the performance of RAID-0 without
 633 scaling operations. The measured array is made up of four disks. Figure 25 plots local
 634 average latencies with and without RAID scaling under the Fin workload, as the time
 635 increases along the x -axis. The average latency without scaling is 1.43 ms, while the
 636 average latency during RAID scaling is 8.01 ms. Figure 26 plots local average latencies
 637 with and without RAID scaling under the TPC-C workload. The average latency
 638 without scaling is 1.38 ms, while the average latency during RAID scaling is 10.39 ms.

639 Obviously, application I/O latencies during RAID scaling are higher than those without
 640 scaling. The reason behind this phenomenon is that application I/Os are continuous

Table II. Comparison of the Improvement of Redistribution Times with Different Disk Sizes

Disk Size	4 GB	8 GB	16 GB
SLAS	6,802s	13,603s	27,206s
FastScale	901s	1,801s	3,598s
Improvement Percentage	86.75%	86.76%	86.77%

641 without long idle periods, either under the Fin workload or under the TPC-C workload.
 642 Interleaving of redistribution I/Os will increase the time of application I/Os waiting
 643 for processing, and the time of disk seeks. The rate-control parameter can be used to
 644 trade off between the redistribution time objective and the response time objective.
 645 In order to obtain acceptable application I/O latencies, one can adjust the rate-control
 646 parameter—the IOPS threshold. In other words, FastScale can accelerate RAID-0 scal-
 647 ing, and at the same time, have an acceptable impact on the existing workload.

648 Running a simulation experiment is time consuming. We set the disk capacity to
 649 4 GB for the online scaling, so as to conduct an experiment in an acceptable time.
 650 We also perform some experiments of offline scaling, where the disk capacity is set
 651 8 GB and 16 GB. As shown in Table II, the redistribution time increases linearly with
 652 the disk size used, no matter which approach is used. However, the percentages of
 653 improvement are consistent with those with the 4 GB capacity.

654 *4.3.2. The Performance after Scaling.* The preceding experiments show that FastScale
 655 improves the scaling efficiency of RAID significantly. One of our concerns is whether
 656 there is a penalty in the performance of the data layout after scaling using FastScale,
 657 compared with the round-robin layout preserved by SLAS.

658 We use the Web workload to measure the performance of the two RAIDs, scaled
 659 from the same RAID using SLAS and FastScale. Each experiment lasts 500 seconds,
 660 and records the latency of each I/O. Based on the issue time, the I/O latency sequence
 661 is evenly divided into 20 sections. Furthermore, we get a local average latency from
 662 each section.

663 First, we compare the performance of the two RAIDs, after one scaling operation
 664 “4+1” using the two scaling approaches. Figure 27 plots local average latencies for the
 665 two RAIDs as the time increases along the x -axis. We find that the performance of the
 666 two RAIDs are very close. With regard to the round-robin RAID, the average latency
 667 is 11.36 ms. For the FastScale RAID, the average latency is 11.37 ms.

668 Second, we compare the performance of the two RAIDs, after two scaling operations
 669 “4+1+1” using the two approaches. Figure 28 plots local average latencies of the two
 670 RAIDs as the time increases along the x -axis. It again reveals approximate equality in
 671 the performance of the two RAIDs. With regard to the round-robin RAID, the average
 672 latency is 11.21 ms. For the FastScale RAID, the average latency is 11.03 ms.

673 Third, we compare the performance of the two RAIDs, after three scaling opera-
 674 tions “4+1+1+1” using the two approaches. Figure 29 plots local average latencies
 675 of the two RAIDs as time increases along the x -axis. It again reveals the approxi-
 676 mate equality in the performance of the two RAIDs. With regard to the round-robin
 677 RAID, the average latency is 11.01 ms. For the FastScale RAID, the average latency is
 678 10.79 ms.

679 Finally, we compare the performance of the two RAIDs, after four scaling opera-
 680 tions “4+1+1+1+1” using the two approaches. Figure 30 plots local average latencies
 681 of the two RAIDs as time increases along the x -axis. It again reveals the approxi-
 682 mate equality in the performance of the two RAIDs. With regard to the round-robin
 683 RAID, the average latency is 10.75 ms. For the FastScale RAID, the average latency is
 684 10.63 ms.

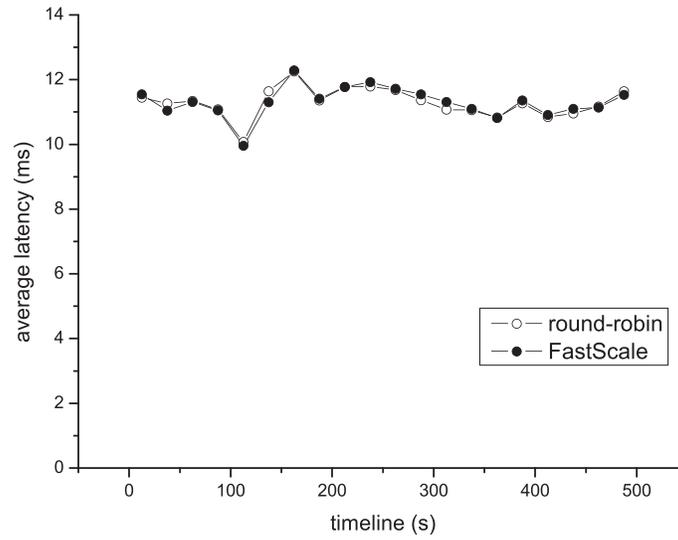


Fig. 27. Performance comparison between FastScale’s layout and round-robin layout under the Web workload after one scaling operation, i.e., “4+1”.

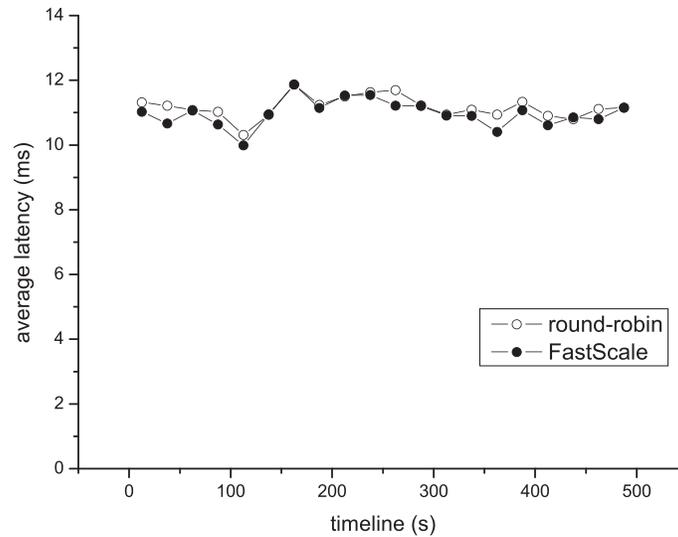


Fig. 28. Performance comparison between FastScale’s layout and round-robin layout under the Web workload after two scaling operations, i.e., “4+1+1”.

685 To summarize, Figure 31 shows a comparison in the response times of the two
 686 RAIDs, scaled from the same RAID using SLAS and FastScale, as the number of scaling
 687 times increases along the x -axis. We can see that the response time of each RAID
 688 decreases as scaling times increase. This is due to an increase of the number of disks
 689 in a RAID, which serve user I/Os simultaneously. The other conclusion that can be
 690 reached is that the two RAIDs, scaled from the same RAID using SLAS and FastScale,
 691 have almost identical performance. In three of the four cases, the FastScale RAID even
 692 performs better than the round-robin RAID.

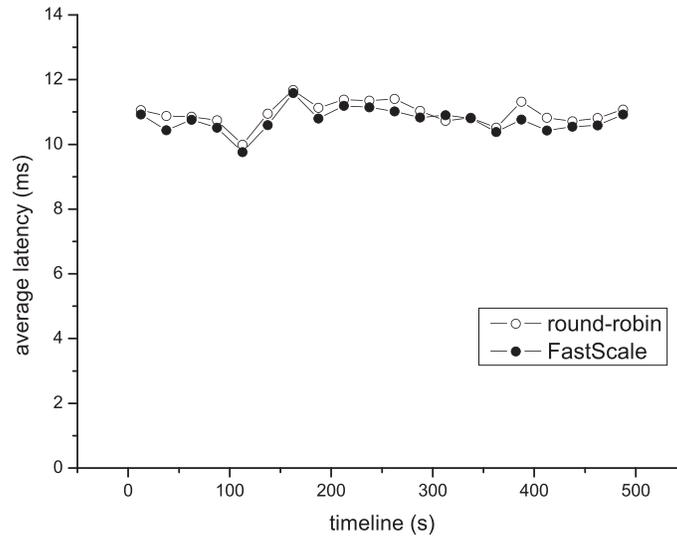


Fig. 29. Performance comparison between FastScale’s layout and round-robin layout under the Web workload after three scaling operations, i.e., “4+1+1+1”.

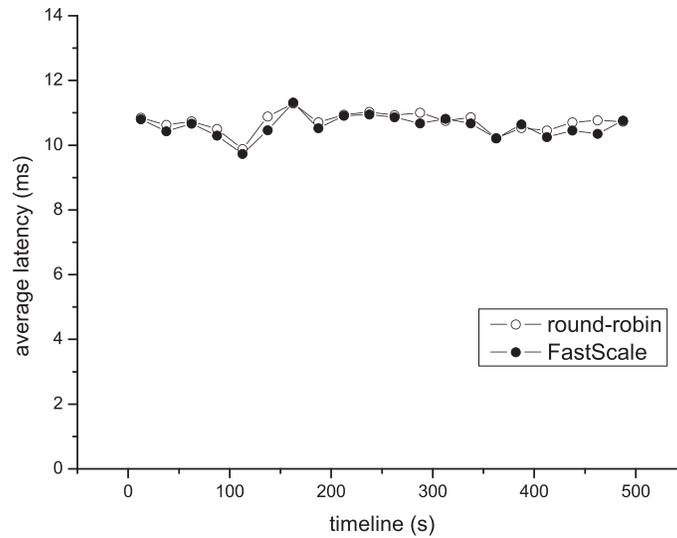


Fig. 30. Performance comparison between FastScale’s layout and round-robin layout under the Web workload after four scaling operations, i.e., “4+1+1+1+1”.

693 5. RELATED WORK

694 In this section, we first examine the existing approaches to scaling deterministic RAID.
 695 Then, we analyze some approaches to scaling randomized RAID.

696 5.1. Scaling Deterministic RAID

697 HP AutoRAID [Wilkes et al. 1996] allows an online capacity expansion. Newly cre-
 698 ated RAID-5 volumes use all of the disks in the system, but previously created RAID-5
 699 volumes continue to use only the original disks. This expansion does not require data

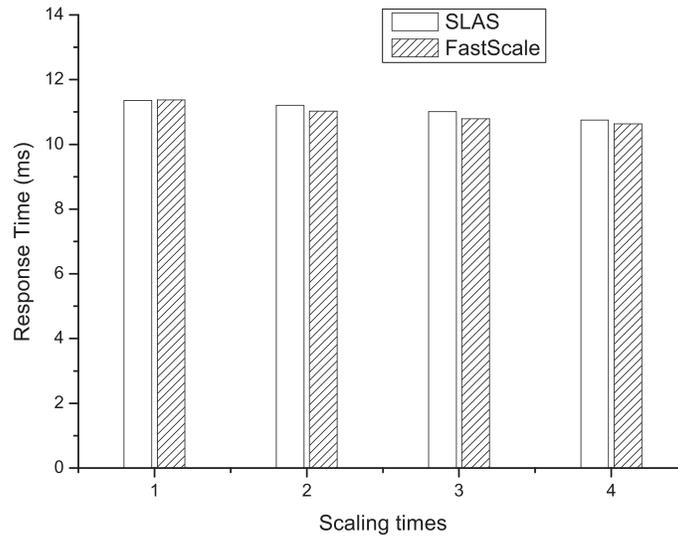


Fig. 31. Comparison of response times of the two RAIDs, scaled from the same RAID using SLAS and FastScale.

700 migration. However, the system cannot add new disks into an existing RAID-5 vol-
 701 ume. The conventional approaches to RAID scaling redistribute data and preserve the
 702 round-robin order after adding disks.

703 Gonzalez and Cortes [2004] proposed a gradual assimilation algorithm (GA) to con-
 704 trol the overhead of scaling a RAID-5 volume. However, GA accesses only one block
 705 via an I/O. Moreover, it writes mapping metadata onto disks immediately after redis-
 706 tributing each stripe. As a result, GA has a large redistribution cost.

707 The reshape toolkit in the Linux MD driver (MD-Reshape) [Brown 2006] writes map-
 708 ping metadata for each fixed-sized data window. However, user requests to the data
 709 window have to queue up until all data blocks within the window are moved. On the
 710 other hand, MD-Reshape issues very small (4KB) I/O operations for data redistribu-
 711 tion. This limits the redistribution performance due to more disk seeks.

712 Zhang et al. [2007] discovered that there is always a reordering window during data
 713 redistribution for round-robin RAID scaling. The data inside the reordering window
 714 can migrate in any order without overwriting any valid data. By leveraging this in-
 715 sight, they proposed the SLAS approach, improving the efficiency of data redistribu-
 716 tion. However, SLAS still requires migrating all data. Therefore, RAID scaling remains
 717 costly.

718 D-GRAID [Sivathanu et al. 2004] restores only live file system data to a hot spare so
 719 as to recover from failures quickly. Likewise, it can accelerate the redistribution pro-
 720 cess if only the live data blocks from the perspective of file systems are redistributed.
 721 However, this requires semantically-smart storage systems. On the contrary, FastScale
 722 is independent of file systems, and it can work with any ordinary disk storage.

723 A patent [Legg 1999] presented a method to eliminate the need to rewrite the origi-
 724 nal data blocks and parity blocks on original disks. However, the method makes all
 725 the parity blocks be either only on original disks or only on new disks. The obvious
 726 distribution nonuniformity of parity blocks will bring a penalty to write performance.

727 Franklin and Wong [2006] presented a RAID scaling method using spare space with
 728 immediate access to new space. First, old data are distributed among the set of data
 729 disk drives and at least one new disk drive while, at the same time, new data are

730 mapped to the spare space. Upon completion of the distribution, new data are copied
731 from the spare space to the set of data disk drives. This is similar to the key idea of
732 WorkOut [Wu et al. 2009]. This kind of method requires spare disks to be available in
733 the RAID.

734 In another patent, Hetzler [2008] presented a method of RAID-5 scaling, called
735 MDM. MDM exchanges some data blocks between original disks and new disks. MDM
736 can perform RAID scaling with reduced data movement. However, it does not increase
737 (just maintains) the data storage efficiency after scaling. The RAID scaling process
738 exploited by FastScale is favored in this regard, because the data storage efficiency is
739 maximized, which many practitioners consider desirable.

740 AdaptiveZ [Gonzalez and Cortes 2007] divides the space of RAID into several adap-
741 tive zones, whose stripe patterns can be customized separately. When new disks are
742 added, AdaptiveZ adds a new zone and redistributes a part of the data at the end of
743 the RAID in the zone. This results in more blocks allocated on old disks than new
744 ones, while redistributing the minimum amount of blocks that the AdaptiveZ algo-
745 rithm requires. Therefore, AdaptiveZ has to increase the size of the migrated zone. In
746 other words, AdaptiveZ is faced with a dilemma between minimal data migration and
747 even data distribution. On the contrary, FastScale combines minimal data migration
748 and uniform data distribution. On the other hand, all the data migrated by AdaptiveZ
749 are logically sequential. On account of spatial locality in I/O workloads, this will still
750 comprise a balanced load.

751 5.2. Scaling Randomized RAID

752 Randomized RAID [Alemany and Thathachar 1997; Brinkmann et al. 2000; Goel et al.
753 2002; Santos et al. 2000] appears to have better scalability. It is now gaining the
754 spotlight in the data placement area. Brinkmann et al. [2000] proposed the cut-and-
755 paste placement strategy that uses randomized allocation strategy to place data across
756 disks. For a disk addition, it cuts off the range $[1/(n + 1), 1/n]$ from given n disks, and
757 pastes them to the newly added $(n + 1)$ th disk. For a disk removal, it uses a reversing
758 operation to move all the blocks in disks that will be removed to the other disks. Also
759 based on random data placement, Seo and Zimmermann [2005] proposed an approach
760 to finding a sequence of disk additions and removals for the disk replacement problem.
761 The goal is to minimize the data migration cost. Both of these approaches assume the
762 existence of a high-quality hash function that assigns all the data blocks in the system
763 to uniformly distributed real numbers with high probability. However, they did not
764 present such a hash function.

765 The SCADDAR algorithm [Goel et al. 2002] uses a pseudo-random function to dis-
766 tribute data blocks randomly across all disks. It keeps track of the locations of data
767 blocks after multiple disk reorganizations and minimizes the amount of data to be
768 moved. Unfortunately, the pseudo-hash function does not preserve the randomness of
769 the data layout after several disk additions or deletions [Seo and Zimmermann 2005].
770 So far, a truly randomized hash function that preserves its randomness after several
771 disk additions or deletions has not been found.

772 The simulation report in Alemany and Thathachar [1997] shows that a single copy
773 of data in random striping may result in some hiccups of the continuous display. To
774 address this issue, one can use data replication [Santos et al. 2000], where a fraction of
775 the data blocks randomly selected are replicated on randomly selected disks. However,
776 this will incur a large overhead.

777 RUSH [Honicky and Miller 2003, 2004] and CRUSH [Weil et al. 2006] are two al-
778 gorithms for online placement and reorganization of replicated data. They are proba-
779 bilistically optimal in distributing data evenly and minimizing data movement when
780 new storage is added to the system. There are three differences between them and

781 FastScale. First, they depend on the existence of a high-quality random function, which
782 is difficult to generate. Second, they are designed for object-based storage systems.
783 They focus on how a data object is mapped to a disk, without considering the data
784 layout of each individual disk. Third, our mapping function needs to be 1-1 and on, but
785 hash functions have collisions and count on some amount of sparseness.

786 6. CONCLUSIONS AND FUTURE WORK

787 This article presents the design of a new approach called FastScale, which accelerates
788 RAID-0 scaling by minimizing data migration. First, with a new and elastic address-
789 ing function, FastScale minimizes the number of data blocks to be migrated without
790 compromising the uniformity of data distribution. Second, FastScale uses access ag-
791 gregation and lazy checkpoint to optimize data migration.

792 Replaying real-system disk I/O traces, we evaluated the performance of FastScale
793 through comparison with an efficient scaling approach called SLAS. The results from
794 detailed experiments show that FastScale can reduce redistribution time by up to
795 86.06% with smaller maximum response time of user I/Os. The experiments also il-
796 lustrate that the performance of the RAID scaled using FastScale is almost identical
797 to, or even better than, that of round-robin RAID.

798 In this article, the factor of data parity is not taken into account. We believe that
799 FastScale provides a good starting point for efficient scaling of RAID-4 and RAID-5
800 arrays. In the future, we will focus on the scaling issues of RAID-4 and RAID-5.

801 ACKNOWLEDGMENTS

802 We would like to thank the three anonymous reviewers for their constructive comments which have helped
803 to improve the quality and presentation of this article.

804 REFERENCES

- 805 Alemany, J. and Thathachar, J. S. 1997. Random striping news on demand servers. Tech. rep. TR-97-02-02,
806 University of Washington.
- 807 Brigham Young University. 2010. TPC-C Postgres 20 iterations. DTB v1.1. Performance Evaluation Labora-
808 tory, Trace distribution center. <http://tds.cs.byu.edu/tds/>.
- 809 Brinkmann, A., Salzwedel, K., and Scheideler, C. 2000. Efficient, distributed data placement strategies for
810 storage area networks. In *Proceedings of the ACM Symposium on Parallel Algorithms and Architectures*.
811 119–128.
- 812 Brown, N. 2006. Online RAID-5 resizing. `drivers/md/raid5.c` in the source code of Linux Kernel 2.6.18.
813 <http://www.kernel.org/>.
- 814 Bucy, J., Schindler, J., Schlosser, S., and Ganger, G. 2008. The DiskSim Simulation Environment Version 4.0
815 Reference Manual. Tech. rep. CMU-PDL-08-101, Carnegie Mellon University.
- 816 Franklin, C. R. and Wong, J. T. 2006. Expansion of RAID subsystems using spare space with immediate
817 access to new space. US Patent 10/033,997.
- 818 Goel, A., Shahabi, C., Yao, S., and Zimmermann, R. 2002. SCADDAR: An efficient randomized technique
819 to reorganize continuous media blocks. In *Proceedings of the 18th International Conference on Data*
820 *Engineering (ICDE)*. 473–482.
- 821 Gonzalez, J. L. and Cortes, T. 2004. Increasing the capacity of RAID5 by online gradual assimilation. In
822 *Proceedings of the International Workshop on Storage Network Architecture and Parallel I/Os (SNAPI)*.
823 17–24.
- 824 Gonzalez, J. L. and Cortes, T. 2007. Adaptive data block placement based on deterministic zones (AdaptiveZ).
825 In *Lecture Notes in Computer Science*, vol. 4804, 1214–1232.
- 826 Hennessy, J. and Patterson, D. 2003. *Computer Architecture: A Quantitative Approach*, 3rd Ed. Morgan
827 Kaufmann Publishers, Inc., San Francisco, CA.
- 828 Hetzler, S. R. 2008. Data storage array scaling method and system with minimal data movement. US Patent
829 20080276057.

- 830 Hitachi. 2001. Hard disk drive specifications Ultrastar 36Z15.
831 [http://www.hitachigst.com/tech/techlib.nsf/techdocs/85256AB8006A31E587256A7800739FEB/\\$file/](http://www.hitachigst.com/tech/techlib.nsf/techdocs/85256AB8006A31E587256A7800739FEB/$file/U36Z15%20sp10.PDF)
832 [U36Z15 sp10.PDF](http://www.hitachigst.com/tech/techlib.nsf/techdocs/85256AB8006A31E587256A7800739FEB/$file/U36Z15%20sp10.PDF). Revision 1.0, April.
- 833 Honicky, R. J. and Miller, E. L. 2003. A fast algorithm for online placement and reorganization of replicated
834 data. In *Proceedings of the 17th International Parallel and Distributed Processing Symposium*.
- 835 Honicky, R. J. and Miller, E. L. 2004. Replication under scalable hashing: A family of algorithms for scal-
836 able decentralized data distribution. In *Proceedings of the 18th International Parallel and Distributed*
837 *Processing Symposium*.
- 838 Kim, C., Kim, G., and Shin, B. 2001. Volume management in SAN environment. In *Proceedings of the 8th*
839 *International Conference on Parallel and Distributed Systems (ICPADS)*. 500–505.
- 840 Legg, C. B. 1999. Method of increasing the storage capacity of a level five RAID disk array by adding, in
841 a single step, a new parity block and N-1 new data blocks which respectively reside in new columns,
842 where N is at least two. US Patent: 6000010, December 1999.
- 843 Muller, K. and Vignaux, T. 2009. SimPy 2.0.1 documentation.
844 <http://simpy.sourceforge.net/SimPyDocs/index.html>.
- 845 Patterson, D. A. 2002. A simple way to estimate the cost of down-time. In *Proceedings of the 16th Large*
846 *Installation Systems Administration Conference (LISA)*. 185–188.
- 847 Patterson, D. A., Gibson, G. A., and Katz, R. H. 1988. A case for redundant arrays of inexpensive disks
848 (RAID). In *Proceedings of the International Conference on Management of Data (SIGMOD)*. 109–116.
- 849 Santos, J. R., Muntz, R. R., and Ribeiro-Neto, B. A. 2000. Comparing random data allocation and data
850 striping in multimedia servers. *ACM SIGMETRICS Perform. Eval. Rev.* 28, 1, 44–55.
- 851 Seo, B. and Zimmermann, R. 2005. Efficient disk replacement and data migration algorithms for large disk
852 subsystems. *ACM Trans. Storage* 1, 3, 316–345.
- 853 Sivathanu, M., Prabhakaran, V., Arpaci-Dusseau, A. C., and Arpaci-Dusseau, R. H. 2004. Improving storage
854 system availability with D-GRAID. In *Proceedings of the 3rd USENIX Conference on File and Storage*
855 *Technologies (FAST)*.
- 856 Storage Performance Council. 2010. <http://www.storageperformance.org/home>.
- 857 UMass Trace Repository. 2007. OLTP Application I/O and Search Engine I/O.
858 <http://traces.cs.umass.edu/index.php/Storage/Storage>.
- 859 Weil, S. A., Brandt, S. A., Miller, E. L., and Maltzahn, C. 2006. CRUSH: Controlled, scalable, decentralized
860 placement of replicated data. In *Proceedings of the International Conference on Supercomputing (SC)*.
- 861 Wilkes, J., Golding, R., Staelin, C., and Sullivan, T. 1996. The HP AutoRAID hierarchical storage system.
862 *ACM Trans. Comput. Syst.* 14, 1, 108–136.
- 863 Wu, S. J., Jiang, H., Feng, D., Tian, L., and Mao, B. 2009. WorkOut: I/O workload outsourcing for boosting the
864 RAID reconstruction performance. In *Proceedings of the 7th USENIX Conference on File and Storage*
865 *Technologies (FAST)*. 239–252.
- 866 Zhang, G. Y., Shu, J. W., Xue, W., and Zheng, W. M. 2007. SLAS: An efficient approach to scaling round-robin
867 striped volumes. *ACM Trans. Storage* 3, 1, 1–39.
- 868 Zheng, W. M. and Zhang, G. Y. 2011. FastScale: Accelerate RAID scaling by minimizing data migration. In
869 *Proceedings of the 9th USENIX Conference on File and Storage Technologies (FAST)*.

870 Received March 2012; revised August 2012; accepted May 2013