

# Reliable and Randomized Data Distribution Strategies for Large Scale Storage Systems

Alberto Miranda\* Sascha Effert† Yangwook Kang‡ Ethan L. Miller‡ Andre Brinkmann† Toni Cortes\*§  
alberto.miranda@bsc.es sascha.effert@upb.de ywkang@cs.ucsc.edu elm@cs.ucsc.edu brinkman@upb.de toni.cortes@bsc.es

\*Barcelona Supercomputing  
Center (BSC)  
Barcelona, Spain

†University of  
Paderborn  
Paderborn, Germany

‡University of  
California  
Santa Cruz, CA, USA

§Technical University of  
Catalonia (UPC)  
Barcelona, Spain

**Abstract**—The ever-growing amount of data requires highly scalable storage solutions. The most flexible approach is to use storage pools that can be expanded and scaled down by adding or removing storage devices. To make this approach usable, it is necessary to provide a solution to locate data items in such a dynamic environment. This paper presents and evaluates the Random Slicing strategy, which incorporates lessons learned from table-based, rule-based, and pseudo-randomized hashing strategies and is able to provide a simple and efficient strategy that scales up to handle exascale data. Random Slicing keeps a small table with information about previous storage system insert and remove operations, drastically reducing the required amount of randomness while delivering a perfect load distribution.

## I. INTRODUCTION

The ever-growing creation of and demand for massive amounts of data requires highly scalable storage solutions. The most flexible approach is to use a pool of storage devices that can be expanded and scaled down as needed by adding new storage devices or removing older ones; this approach necessitates a scalable solution for locating data items in such a dynamic environment.

Table-based strategies can provide an optimal mapping between data blocks and storage systems, but obviously do not scale to large systems because tables grow linearly in the number of data blocks. Rule-based methods, on the other hand, run into fragmentation problems, so defragmentation must be performed periodically to preserve scalability.

Hashing-based strategies use a compact function  $h$  in order to map balls with unique identifiers out of some large universe  $U$  into a set of bins called  $S$  so that the balls are evenly distributed among the bins. In our case, balls are data items and bins are storage devices. Given a static set of devices, it is possible to construct a hash function so that every device gets a fair share of the data load. However, standard hashing techniques do not adapt well to a changing set of devices.

Consider, for example, the hash function  $h(x) = (a \cdot x + b) \bmod n$ , where  $S = \{0, \dots, n - 1\}$  represents the set of

This work was partially supported by the Spanish Ministry of Science and Technology under the TIN2007-60625 grant, the Catalan Government under the 2009-SGR-980 grant, the EU Marie Curie Initial Training Network SCALUS under grant agreement no. 238808, the National Science Foundation under grants CCF-0937938 and IIP-0934401, and by the industrial sponsors of the Storage Systems Research Center at the University of California, Santa Cruz.

storage devices. If a new device is added, we are left with two choices: either replace  $n$  by  $n + 1$ , which would require virtually all the data to be relocated; or add additional rules to  $h(x)$  to force a certain set of data blocks to be relocated on the new device in order to get back to a fair distribution, which, in the long run, destroys the compactness of the hashing scheme.

Pseudo-randomized hashing schemes that can adapt to a changing set of devices have been proposed and theoretically analyzed. The most popular is probably Consistent Hashing [17], which is able to evenly distribute single copies of each data block among a set of storage devices and to adapt to a changing number of disks. We will show that these pure randomized data distribution strategies have, despite their theoretical perfectness, serious drawbacks when used in very large systems.

Besides adaptivity and fairness, redundancy is important as well. Storing just a single copy of a data item in real systems is dangerous because, if a storage device fails, all of the blocks stored in it are lost. It has been shown that simple extensions of standard randomized data distribution strategies to store more than a single data copy are not always capacity efficient [5].

The main contributions of this paper are:

- **First comparison of different hashing-based data distribution strategies that are able to replicate data in a heterogeneous and dynamic environment.** This comparison shows the strengths and drawbacks of the different strategies as well as their constraints. Such comparison is novel because hashing-based data distribution strategies have been mostly analytically discussed, with only a few implementations available, and in the context of peer-to-peer networks with limited concern for the fairness of the data distribution [23]. Only a few of these strategies have been implemented in storage systems, where limited fairness immediately leads to a strong increase in costs [6][24].
- The introduction of **Random Slicing, which overcomes the drawbacks of randomized data distribution strategies** by incorporating lessons learned from table-based, rule-based and pseudo-randomized hashing strategies. Random Slicing keeps a small table with information about previous storage system insertions and removals. This table helps to drastically reduce the required amount

of randomness in the system and thus reduces the amount of necessary main memory by orders of magnitude.

It is important to note that all randomized strategies map (virtual) addresses to a set of disks, but do not define the placement of the corresponding block on the disk surface. This placement on the block devices has to be resolved by additional software running on the disk itself. Therefore, we will assume inside the remainder of the paper that the presented strategies work in an environment that uses object-based storage. Unlike conventional block-based hard drives, object-based storage devices (OSDs) manage disk block allocation internally, exposing an interface that allows others to read and write to variably-sized, arbitrarily-named objects [2][12].

### A. The Model

Our research is based on an extension of the standard “balls into bins” model [16][19]. Let  $\{0, \dots, M - 1\}$  be the set of all identifiers for the balls and  $\{0, \dots, N - 1\}$  be the set of all identifiers for the bins. Suppose that the current number of balls in the system is  $m \leq M$  and that the current number of bins in the system is  $n \leq N$ . We will often assume for simplicity that the balls and bins are numbered in a consecutive way starting with 0, but any numbering that gives unique numbers to each ball and bin would work for our strategies.

Suppose that bin  $i$  can store up to  $b_i$  (copies of) balls. Then we define its relative capacity as  $c_i = b_i / \sum_{j=0}^{n-1} b_j$ . We require that, for every ball,  $k$  copies must be stored on different bins for some fixed  $k$ . In this case, a trivial upper bound for the number of balls the system can store while preserving fairness and redundancy is  $\sum_{j=0}^{n-1} b_j / k$ , but it can be much less than that in certain cases. We term the  $k$  copies of a ball a *redundancy group*.

Placement schemes for storing redundant information can be compared based on the following criteria (see also [8]):

- **Capacity Efficiency and Fairness:** A scheme is called *capacity efficient* if it allows us to store a near-maximum number of data blocks. We will see in the following that the fairness property is closely related to capacity efficiency, where *fairness* describes the property that the number of balls **and** requests received by a bin are proportional to its capacity.
- **Time Efficiency:** A scheme is called *time efficient* if it allows a fast computation of the position of any copy of a data block without the need to refer to centralized tables. Schemes often use smaller tables that are distributed to each node that must locate blocks.
- **Compactness:** We call a scheme *compact* if the amount of information the scheme requires to compute the position of any copy of a data block is small (in particular, it should only depend on  $n$ —the number of bins).
- **Adaptivity:** We call a scheme *adaptive* if it only redistributes a near-minimum amount of copies when new storage is added in order to get back into a state of fairness. We therefore compare the different strategies in Section V with the minimum amount of movements, which is required to keep the fairness property.

Our goal is to find strategies that perform well under all of these criteria.

### B. Previous Results

Data reliability and support for scalability as well as the dynamic addition and removal of storage systems is one of the most important issues in designing storage environments. Nevertheless, up to now only a limited number of strategies has been published for which it has formally been shown that they can perform well under these requirements.

Data reliability is achieved by using RAID encoding schemes, which divide data blocks into specially encoded sub-blocks that are placed on different disks to make sure that a certain number of disk failures can be tolerated without losing any information [20]. RAID encoding schemes are normally implemented by striping data blocks according to a pre-calculated pattern across all the available storage devices. Even though deterministic extensions for the support of heterogeneous disks have been developed [10][13], adapting the placement to a changing number of disks is cumbersome under RAID as all of the data may have to be reorganized.

In the following, we just focus on data placement strategies that are able to cope with dynamic changes of the capacities or the set of storage devices in the system. Karger, *et al.* present an adaptive hashing strategy for homogeneous settings that satisfies fairness and is 1-competitive w.r.t. adaptivity [17]. In addition, the computation of the position of a ball takes only an expected number of  $O(1)$  steps. However, their data structures need at least  $n \log^2 n$  bits to ensure a good data distribution.

Brinkmann, *et al.* presented the cut-and-paste strategy as alternative placement strategy for uniform capacities [7]. Their scheme requires  $O(n \log n)$  bits and  $O(\log n)$  steps to evaluate the position of a ball. Furthermore, it keeps the deviation from a fair distribution of the balls extremely small with high probability. Interestingly, the theoretical analysis of this strategy has been experimentally re-evaluated in a recent paper by Zheng *et al.* [26].

Sanders considers the case that bins fail and suggests to use a set of forwarding hash functions  $h_1, h_2, \dots, h_k$ , where at the time  $h_i$  is set up, only bins that are intact at that time are included in its range [21].

Adaptive data placement schemes that are able to cope with arbitrary heterogeneous capacities have been introduced in [8]. The presented strategies Share and Sieve are compact, fair, and (amortized)  $(1 + \epsilon)$ -competitive for arbitrary changes from one capacity distribution to another, where  $\epsilon > 0$  can be made arbitrarily small. Other data placement schemes for heterogeneous capacities are based on geometrical constructions [22]; the linear method used combines the standard consistent hashing approach [17] with a linear weighted distance measure.

All previously mentioned work is only applicable for environments where no replication is required. Certainly, it is easy to come up with proper extensions of the schemes so that no two copies of a ball are placed in the same bin. A simple approach feasible for all randomized strategies to replicate a ball  $k$  times is to perform the experiment  $k$  times and to

remove after each experiment the selected bin. Nevertheless, it has been shown that the fairness condition cannot be guaranteed for these simple strategies and that capacity will be wasted [5]. This paper will also evaluate the influence of this capacity wasting in realistic settings.

The first methods with dedicated support for replication were proposed by Honicky and Miller [14][15]. *RUSH* (Replication Under Scalable Hashing) maps replicated objects to a scalable collection of storage servers according to user-specified server weighting. When the number of servers changes, RUSH tries to redistribute as few objects as possible to restore a balanced data distribution while ensuring that no two replicas of an object are ever placed on the same server.

*CRUSH* is derived from RUSH, and supports different hierarchy levels that provide the administrator finer control over the data placement in the storage environment [25]. The algorithm accommodates a wide variety of data replication and reliability mechanisms and distributes data in terms of user-defined policies.

Amazon’s Dynamo [11] uses a variant of Consistent Hashing with support for replication where each node is assigned multiple tokens chosen at random that are used to partition the hash space. This variant has given good results concerning performance and fairness, though the authors claim that it might have problems scaling up to thousands of nodes.

Brinkmann *et al.* have shown that a huge class of placement strategies cannot preserve fairness and redundancy at the same time and have presented a placement strategy for an arbitrary fixed number  $k$  of copies for each data block, which is able to run in  $O(k)$ . The strategies have a competitiveness of  $\log n$  for the number of replacements in case of a change of the infrastructure [5]. This competitiveness has been reduced to  $O(1)$  by breaking the heterogeneity of the storage systems [4]. Besides the strategies presented inside this paper, it is worth mentioning the Spread strategy, which has similar properties to those of Redundant Share [18].

## II. RANDOMIZED DATA DISTRIBUTION

We present in this section a short description of the applied data distribution strategies. We start with Consistent Hashing and Share, which can, in their original form, only be applied for  $k = 1$  and therefore lack support for any redundancy strategy. Both strategies are used as sub-strategies inside some of the investigated data distribution strategies. Besides their usage as sub-strategies, we will also present a simple replication strategy, which can be based on any of these simple strategies. Afterwards, we present Redundant Share and RUSH, which directly support data replication.

### A. Consistent Hashing

We start with the description of the Consistent Hashing strategy, which solves the problem of (re-)distributing data items in homogeneous systems [17]. In Consistent Hashing, both data blocks and storage devices are hashed to random points in a  $[0, 1)$ -interval, and the storage device closest to a data block in this space is responsible for that data

block. Consistent Hashing ensures that adding or removing a storage device only requires a near minimal amount of data replacements to get back to an even distribution of the load. However, the consistent hashing technique cannot be applied well if the storage devices can have arbitrary non-uniform capacities since in this case the load distribution has to be adapted to the capacity distribution of the devices. The memory consumption of Consistent Hashing heavily depends on the required fairness. Using only a single point for each storage device leads to a load deviation of  $n \cdot \log n$  between the least and heaviest loaded storage devices. Instead it is necessary to use  $\log n$  virtual devices to simulate each physical device, respectively to throw  $\log n$  points for each device to achieve a constant load deviation.

### B. Share-strategy

Share supports heterogeneous environments by introducing a two stage process [8]. In the first stage, the strategy randomly maps one interval for each storage system to the  $[0, 1)$ -interval. The length of these intervals is proportional to the size of the corresponding storage systems and some stretch factor  $s$  and can cover the  $[0, 1)$ -interval many times. In this case, the interval is represented by several virtual intervals. The data items are also randomly mapped to a point in the  $[0, 1)$ -interval. Share now uses an adaptive strategy for homogeneous storage systems, like Consistent Hashing, to get the responsible storage systems from all storage systems for which the corresponding interval includes this point.

The analysis of the Share-strategy shows that it is sufficient to have a stretch factor  $s = O(\log N)$  to ensure correct functioning and that Share can be implemented in expected time  $O(1)$  using a space of  $O(s \cdot k \cdot (n + 1/\delta))$  words (without considering the hash functions), where  $\delta$  characterizes the required fairness. Share has an amortized competitive ratio of at most  $1 + \epsilon$  for any  $\epsilon > 0$ . Nevertheless, we will show that, similar to Consistent Hashing, the memory consumption heavily depends on the expected fairness.

### C. Trivial data replication

Consistent Hashing and Share are, in their original setting, unable to support data replication or erasure codes, since it is always possible that multiple strips belonging to the same stripe set are mapped to the same storage system and that data recovery in case of failures becomes impossible. Nevertheless, it is easy to imagine strategies to overcome this drawback and to support replication strategies by, e.g., simply removing all previously selected storage systems for the next random experiment for a stripe set. Another approach, used inside the experiments in this paper, is to simply perform as many experiments as are necessary to get enough independent storage systems for the stripe set. It has been shown that this trivial approach wastes some capacity [5], but we will show in this paper that this amount can often be neglected.

### D. Redundant Share

Redundant Share has been developed to support the replication of data in heterogeneous environments. The strategy

orders the bins according to their weights  $c_i$  and sequentially iterates over the bins [5]. The basic idea is that the weights are calculated in a way that ensures perfect fairness for the first copy and to use a recursive descent to select additional copies. Therefore, the strategy needs  $O(n)$  rounds for each selection process. The algorithm is  $\log n$ -competitive concerning the number of replacements if storage systems enter or leave the system. The authors of the original strategy have also presented extensions of Redundant Share, which are  $O(1)$ -competitive concerning the number of replacements [4] as well as strategies, which have  $O(k)$ -runtime. Both strategies rely on Share and we will discuss in the evaluation section, why they are therefore not feasible in realistic settings.

### E. RUSH

The RUSH algorithms all proceed in two stages, first identifying the appropriate cluster in which to place an object, and then identifying the disk within a cluster. Within a cluster, replicas assigned to the cluster are mapped to disks using prime number arithmetic that guarantees that no two replicas of a single object can be mapped to the same disk. Selection of clusters is a bit more complex, and differs between the three RUSH variants:  $RUSH_P$ ,  $RUSH_R$ , and  $RUSH_T$ .  $RUSH_P$  considers clusters in the reverse of the order they were added, and determines whether an object would have been moved to the cluster when it was added; if so, the search terminates and the object is placed.  $RUSH_R$  works in a similar way, but it determines the number of objects in each cluster simultaneously, rather than requiring a draw for each object.  $RUSH_T$  improves the scalability of the system by descending a tree to assign objects to clusters; this reduces computation time to  $\log c$ , where  $c$  is the number of clusters added.

Given that  $RUSH_R$  and  $RUSH_T$  outperform  $RUSH_P$ , and that both showed similar properties during our analysis, we will only include results for  $RUSH_R$  for the sake of brevity.

## III. RANDOM SLICING

Random Slicing is designed to be fair and efficient both in homogeneous and heterogeneous environments and to adapt gracefully to changes in the number of bins. Suppose that we have a random function  $h : \{1, \dots, M\} \rightarrow [0, 1)$  that maps balls uniformly at random to real numbers in the interval  $[0, 1)$ . Also, suppose that the relative capacities for the  $n$  given bins are  $(c_0, \dots, c_{n-1}) \in [0, 1)^n$  and that  $\sum_{i=0}^{n-1} c_i = 1$ .

The strategy works by dividing the  $[0, 1)$  range into intervals and assigning them to the bins currently in the system. Notice that the intervals created do not overlap and completely cover the  $[0, 1)$  range. Also note that bin  $i$  can be responsible for several non-contiguous intervals  $P_i = (I_0, \dots, I_k)$ , where  $k < n$ , which will form the *partition* of that bin. To ensure fairness, Random Slicing will always enforce that  $\sum_{j=0}^{k-1} |I_j| = c_i$ .

In an initial phase, i.e. when the first set of bins enters the system, each bin  $i$  is given only one interval of length  $c_i$ , since this suffices to maintain fairness. Whenever new bins enter the system, however, relative capacities for old bins change due to the increased overall capacity. To maintain fairness, Random

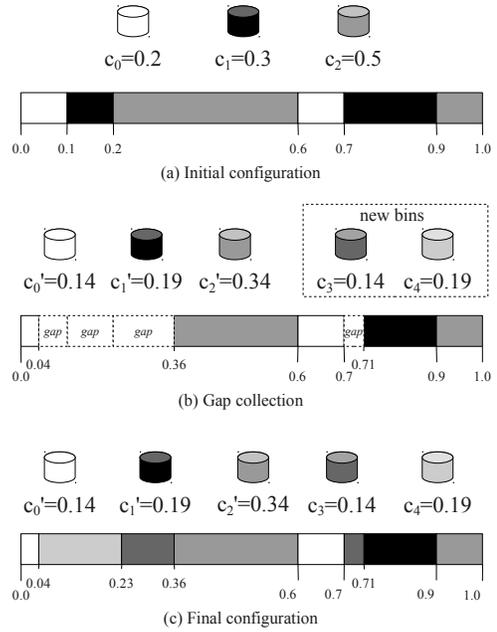


Fig. 1. Random Slicing's interval reorganization when adding new bins.

Slicing shrinks existing partitions by splitting the intervals that compose them until their new relative capacities are reached. The new intervals generated are used to create partitions for the new bins.

Algorithm 1 shows the mechanism used to reduce the interval length of partitions in detail. First, the algorithm computes by how much partitions should be reduced in order to keep the fairness of the distribution. Since the global capacity has increased, each partition  $P_i$  must be reduced by  $r_i = c_i - c'_i$ , where  $c'_i$  corresponds to the new relative capacity of bin  $i$ .

Partitions become smaller by releasing or splitting some of their intervals, thus generating *gaps*, which can be used for new intervals. Notice, however, that the strategy's memory consumption directly depends on the number of intervals used and, therefore, the number of splits made in each addition phase can affect scalability. For this reason, the algorithm tries to collect as many complete intervals as possible and will only split an existing interval as a last resort. Furthermore, when splitting an interval is the only option, the algorithm tries to expand any adjacent gap instead of creating a new one.

The partition lengths for the old bins already represent the corresponding relative capacities. It is only necessary to use these gaps to create new partitions for the newly added bins. The strategy proceeds by greedily allocating the largest partitions to the largest gaps available in order to reduce the number of new intervals even more, which ends the process.

An example of this reorganization is shown in Figure 1, where two new bins  $B_3$  and  $B_4$ , representing a 50% capacity increase, are added to the bins  $B_0$ ,  $B_1$ , and  $B_2$ . Figure 1(a) shows the initial configuration and the relative capacities for the initial bins. Figure 1(b) shows that the partition of  $B_0$  must be reduced by 0.06, the partition of  $B_1$  by 0.11, and

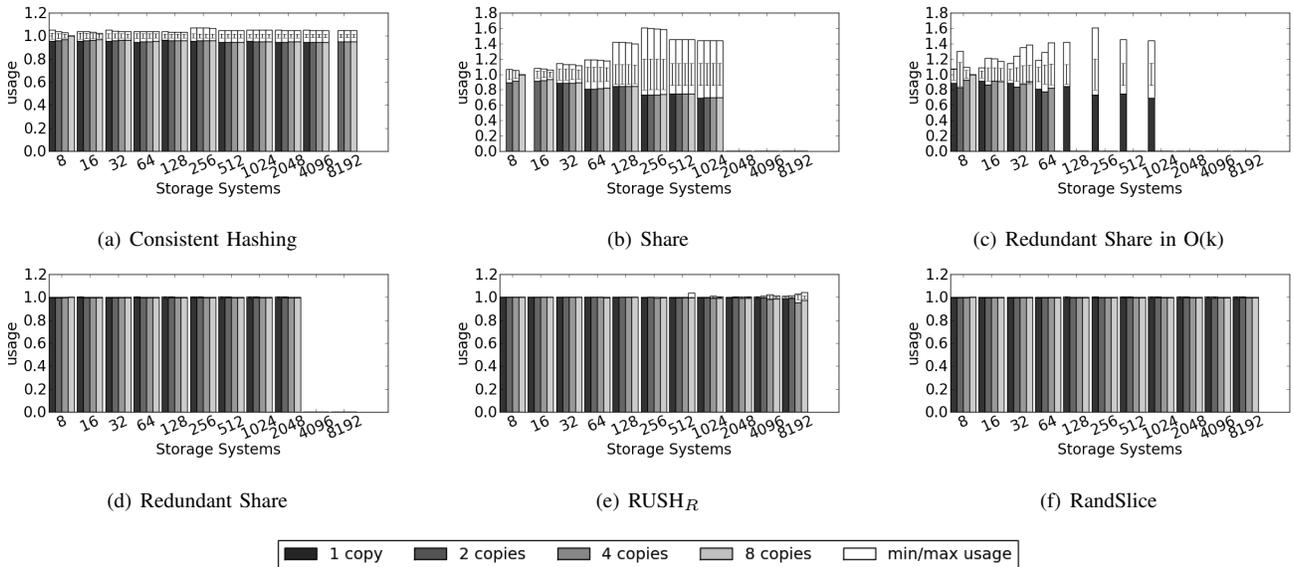


Fig. 2. Fairness of the data distribution strategies for homogeneous disks.

the one of  $B_2$  by 0.16, whereas two new partitions with a size of 0.14 and 0.19 must be created for  $B_3$  and  $B_4$ . The interval  $[0.1, 0.2] \in B_1$  can be completely cannibalized, whereas the intervals  $[0.0, 0.1] \in B_0$ ,  $[0.2, 0.6] \in B_2$  and  $[0.7, 0.9] \in B_1$  are split while trying to maximize gap lengths. Figure 1(c) shows that the partition for  $B_3$  is composed of intervals  $[0.23, 0.36)$  and  $[0.7, 0.71)$ , while the partition for  $B_4$  consists only of interval  $[0.04, 0.23)$ .

When all partitions are created, the location of a ball  $b$  can be determined by calculating  $x = h(b)$  and retrieving the bin associated with it. Notice that some balls will change partition after the reorganization, but as partitions always match their ideal capacity, only a near minimal amount of balls will need to be reallocated. Furthermore, if  $h(b)$  is uniform enough and the number of balls in the system significantly larger than the number of intervals (both conditions easily feasible), the fairness of the strategy is guaranteed.

#### IV. METHODOLOGY

Most previous evaluations of data distribution strategies are based on an analytical investigation of their properties. In contrast, we will use a simulation environment to examine the real-world properties of the investigated protocols. All data distribution strategies will be evaluated in the same environment, which scales from few storage systems up to thousands of devices. The simulation environment has been developed by the authors of this papers and has been made available online<sup>1</sup>. The collection also include the parameter settings for the individual experiments.

We distinguish between homogeneous and heterogeneous settings and also between static and dynamic environments.

Typical storage systems do not work on the individual hard disk level, but are handled based on a granularity of shelves.

We assume that each storage node (also called storage system in the following) consists of 16 hard disks (plus potential disks to add additional intra-shelf redundancy).

We assume that each storage systems in the homogeneous, static setting can hold up to  $k \cdot 500,000$  data items, where  $k$  is the number of copies of each block. Assuming a hard disk capacity of 1 TByte and putting 16 hard disks in each shelf means that each data item has a size of 2 MByte. The number of placed data items is  $k \cdot 250,000$  times the number of storage systems. In all cases, we compare the fairness, the memory consumption, as well as the performance of the different strategies for a different number of storage systems.

The heterogeneous setting assumes that we have 128 storage systems in the beginning and we add in each step 128 systems, which have  $3/2$  times the size of the previously added system. We are placing again half the number of items, which saturates all disks.

For each of the homogeneous and heterogeneous tests, we also count the number of data items, which have to be moved in case we are adding disks, so that the data distribution delivers the correct location for a data item after the redistribution phase. The number of moved items has to be as small as possible to support dynamic environments, as the systems typically tend to a slower performance during the reconfiguration process.

The dynamic behavior can be different if the order of the  $k$  copies is important, e.g. in case of parity RAID, Reed-Solomon codes, or EvenOdd-Codes, or if this order can be neglected in case of pure replication strategies [20][3][9].

#### V. EVALUATION

The following section evaluates the impact of the different distribution strategies on the data distribution quality, the memory consumption of the different strategies, their adaptivity and performance. All graphs presented in the section contain four

<sup>1</sup><http://dadisi.sourceforge.net>

---

**Algorithm 1** Gap Collection in Random Slicing

---

**Input:**  $\{b_0, \dots, b_{n-1}\}, \{b_n, \dots, b_{p-1}\}, \{I_0, \dots, I_{q-1}\}$ **Output:**  $gaps : \{G_0, \dots, G_{m-1}\}$ **Require:**  $(p > n) \wedge (q \geq n)$ 

```
1:  $\forall i \in \{0, \dots, p-1\} : c'_i \leftarrow b_i / \sum_{j=1}^{p-1} b_j$ 
2:  $\forall i \in \{0, \dots, n-1\} : r_i \leftarrow c_i - c'_i$ 
3:  $gaps \leftarrow \{\}$ 
4: for  $i \leftarrow 0$  to  $q-1$  do
5:    $j \leftarrow$  get bin assigned to  $I_i$ 
6:    $G \leftarrow$  get last gap from  $gaps$ 
7:   if  $r_j > 0$  then
8:     if  $length(I_i) < r_j$  then
9:       if  $adjacent(G, I_i)$  then
10:         $G \leftarrow G + length(I_i)$ 
11:      else
12:         $gaps \leftarrow gaps + I_i$ 
13:      end if
14:       $r_i \leftarrow r_i - length(I_i)$ 
15:      if last interval was assimilated completely then
16:         $cut\_interval\_end \leftarrow false$ 
17:      end if
18:    else
19:      if  $adjacent(G, I_i)$  then
20:         $G \leftarrow G + length(I_i)$ 
21:      else
22:        if  $cut\_interval\_end$  then
23:           $gaps \leftarrow gaps + \{I_i.end - r_j, I_i.end\}$ 
24:        else
25:           $gaps \leftarrow gaps + \{I_i.start, I_i.start + r_j\}$ 
26:        end if
27:      end if
28:       $r_i \leftarrow r_i - r_j$ 
29:       $cut\_interval\_end \leftarrow \neg cut\_interval\_end$ 
30:    end if
31:  end if
32: end for
33: return  $gaps$ 
```

---

bars for each number of storage systems, which represent the experimental results for one, two, four, and eight copies (please see Figure 2 for the color codes in the legend). The white boxes in each bar represent the range of results, e.g., between the minimum and the maximum usage. Also, the white boxes include the standard deviation for the experiments. Small or non-existing white boxes indicate a very small deviation between the different experiments. Sometimes we print less relevant or obvious results in smaller boxes to stay inside the page limit.

#### A. Fairness

The first simulations evaluate the fairness of the strategies for different sets of homogeneous disks, ranging from 8 storage systems up to 8192 storage systems (see Figure 2).

Consistent Hashing has been developed to evenly distribute one copy over a set of homogeneous disks of the same

size. Figure 2(a) shows that the strategy is able to fulfill these demands for the test case, in which all disks have the same size. The difference between the maximum and the average usage is always below 7% and the difference between the minimum and average usage is always below 6%. The deviation is nearly independent from the number of copies as well as from the number of disks in the system, so that the strategy can be reasonably well applied. We have thrown  $400 \cdot \log n$  points for each storage system (please see Section II-A for the meaning of points in Consistent Hashing).

The fairness of Consistent Hashing can be improved by throwing more points for each storage system (see Figure 3 for an evaluation with 64 storage systems). The evaluation shows that initial quality improvements can be achieved with very few additional points, while further small improvements require a high number of extra points per storage system.  $400 \cdot \log n$  points are 2400 points for 64 storage systems, meaning that we are already using a high number of points, where further quality improvement becomes very costly.

Share has been developed to overcome the drawbacks of Consistent Hashing for heterogeneous disks. Its main idea is to (randomly) partition the disks into intervals and assign a set of disks to each interval. Inside an interval, each disk is treated as homogeneous and strategies like Consistent Hashing can be applied to finally distribute the data items.

The basic idea implies that Share has to compute and keep the data structures for each interval. 1,000 disks lead to a maximum of 2,000 intervals, implying 2,000 times the memory consumption of the applied uniform strategy. On the other hand, the number of disks inside each interval is smaller than  $n$ , which is the number of disks in the environment. The analysis of Share shows that on average  $c \cdot \log n$  disks participate in each interval (see Section II-B, without loss of generality we will neglect the additional  $\frac{1}{\delta}$  to keep the argumentation simple). Applying Consistent Hashing as homogeneous strategy therefore leads to a memory consumption, which is in  $O(n \cdot \log^2 n \cdot \log^2(\log n))$  and therefore only by a factor of  $\log^2(\log n)$  bigger than the memory consumption of Consistent Hashing.

Unfortunately, it is not possible to neglect the constants in a real implementation. Figure 2(b) shows the fairness of Share for a stretch factor of  $3 \cdot \log n$ , which shows huge deviations even for homogeneous disks. A deeper analysis of the Chernoff-bounds used in [8] shows that it would have been necessary to have a stretch factor of 2,146 to keep fairness

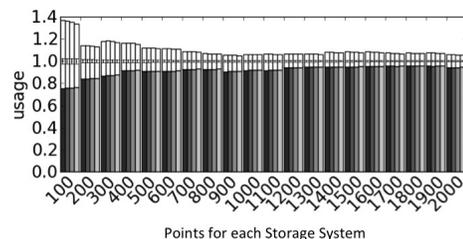


Fig. 3. Influence of point number on Consistent Hashing.

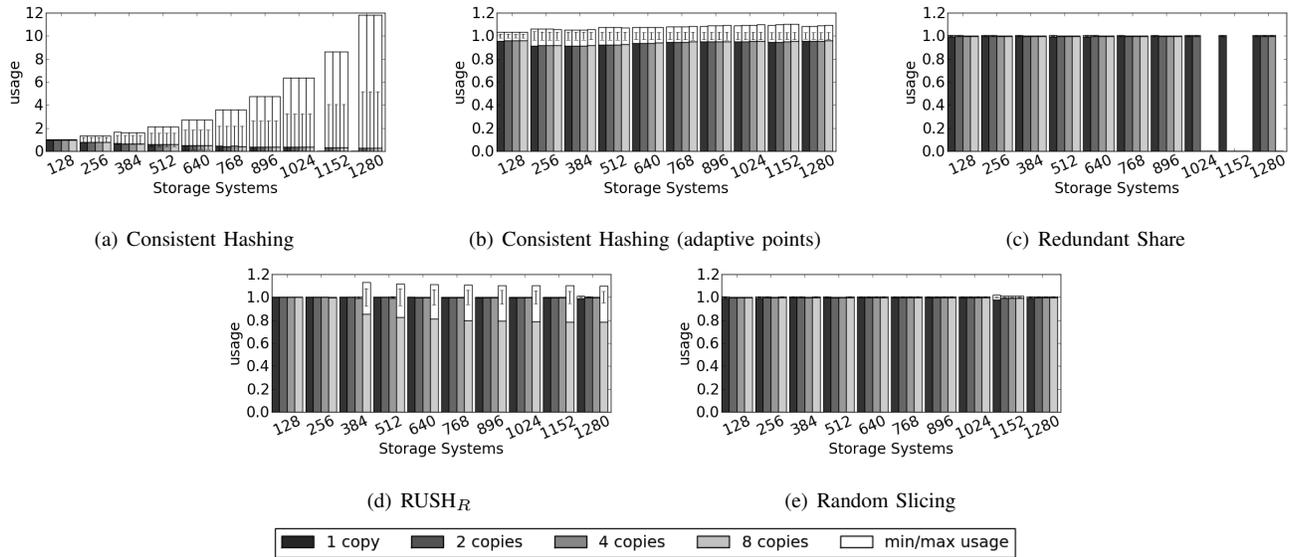


Fig. 5. Fairness of the data distribution strategies for heterogeneous disks.

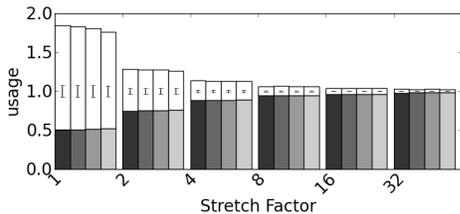


Fig. 4. Fairness of Share depending on stretch factor.

in the same order as the fairness achieved with Consistent Hashing, which is infeasible in scale-out environments.

Simulations including different stretch factors for 64 storage systems for Share are shown in Figure 4, where the x-axis depicts the stretch factor divided by  $\ln n$ . The fairness can be significantly improved by increasing the stretch factor. Unfortunately, a stretch factor of 32 already requires in our simulation environment more than 50 GByte main memory for 64 storage systems, making Share impractical in bigger environments. In the following, we will therefore skip this strategy in our evaluations.

Redundant Share uses precomputed intervals for each disk and therefore does not rely too much on randomization properties. The intervals exactly represent the share of each disk on the total disk capacity, leading to a very even distribution of the data items (see Figure 2(d)). The drawback of this version of Redundant Share is that it has linear runtime, possibly leading to high delays in case of huge environments. Brinkmann et al. have presented enhancements, which enable Redundant Share to have a runtime in  $O(k)$ , where  $k$  is the number of copies [4]. Redundant Share in  $O(k)$  requires a huge number of Share instances as sub-routines, making it impractical to support a huge number of disks and a good fairness at the same time. Figure 2(c) shows that it is even difficult to support

multiple copies for more than 64 disks, even if the required fairness is low, as 64 GByte main memory have not been sufficient to calculate these distributions. Therefore, we will also neglect Redundant Share with runtime in  $O(k)$  in the following measurements.

$RUSH_R$  places objects almost ideally according to the appropriate weights, though it begins to degrade as the number of disks grows (see Figure 2(e)). We believe this happens due to small variations in the probabilistic distribution, which build up for higher numbers of storage systems.

In Random Slicing, precomputed partitions are used to represent a disk’s share of the total system capacity, in a similar way to Redundant Share’s use of intervals. This property, in addition to the hash function used, enforces an almost optimal distribution of the data items, as shown in Figure 2(f).

The fairness of the different strategies for a set of heterogeneous storage systems is depicted in Figure 5. As described in Section IV, we start with 128 storage systems and add every time 128 additional systems having  $3/2$ -times the capacity of the previously added.

The fairness of Consistent Hashing in its original version is obviously very poor (see Figure 5(a)). Assigning the same number of points in the  $[0, 1]$ -interval for each storage system, independent of its size, leads to huge variations. Simply adapting the number of points based on the capacities leads to much better deviations (see Figure 5(b)). The difference between the maximum, respectively minimum and the average usage is around 10% and increases slightly with the number of copies. In the following, we will always use Consistent Hashing with an adaptive number of copies, depending on the capacities of the storage systems.

Both Redundant Share and Random Slicing show again a nearly perfect distribution of data items over the storage systems, due to their precise modeling of disk capacities and the uniformity of the distribution functions (see Figures 5(c)

and 5(e), respectively).

Figure 5(d) shows that  $RUSH_R$  does a good distribution job for 1, 2, and 4 copies but seems to degrade with 8 copies showing important deviations from the optimal distribution.

### B. Memory consumption and compute time

The memory consumption as well as the performance of the different data distribution strategies have a strong impact on the applicability of the different strategies. We assume that scale-out storage systems mostly occur in combination with huge cluster environments, where the different cores of a cluster node can share the necessary data structures for storage management. Assuming memory capacities of 192 GByte per node in 2015 [1], we do not want to waste more than 10% or approximately 20 GByte of this capacity for the metadata information of the underlying storage system. Furthermore, we assume access latencies of 5 ms for a magnetic storage system and access latencies of 50  $\mu s$  for a solid state disks. These access latencies set an upper limit on the time allowed for calculating the translation from a virtual address to a physical storage system.

The bars in the graphs of Figure 6 represent the average allocated memory, the white bars on top the peak consumption of virtual memory over the different tests. The points in that figure represent the average time required for a single request. These latencies include confidence intervals.

The memory consumption of Consistent Hashing only depends on the number and kind of disks in the system, while the number of copies  $k$  has no influence on it (see Figure 6(a)). We are throwing  $400 \cdot \log n$  points for the smallest disk, the number of points for bigger disks grows proportional to their capacity, which is necessary to keep fairness in heterogeneous environments. Using 1,280 heterogeneous storage systems requires a memory capacity of nearly 9 GByte, which is still below our limit of 20 GByte.

The time to calculate the location of data item only depends on the number of copies, as Consistent Hashing is implemented as a  $O(1)$ -strategy for a single copy. Therefore, it is possible to use Consistent Hashing in scale-out environments, which are based on solid state drives, as the average latency for the calculation of a single data item stays below 10  $\mu s$ .

Redundant Share has very good properties concerning memory usage, but the computation time grows linearly in the number of storage systems. Even the calculation of a single item for 128 storage systems takes 145  $\mu s$ . Using 8 copies increases the average access time for all copies to 258  $\mu s$ , which is 50  $\mu s$  for each copy, making it suitable for mid sized environments, which are based on SSDs. Increasing the environment to 1280 storage systems raised the calculation time almost linearly for a single copy to 669  $\mu s$ , which is reasonable in magnetic disk based environments.

$RUSH_R$  shows good results both in memory consumption and in computation time (see Figure 6(c)). The reduced memory consumption is explained because the strategy does not need a great deal of in-memory structures in order to maintain the information about clusters and storage nodes. Lookup

times depend only on the number of clusters in the system, which can be kept comparatively small for large systems.

Random Slicing shows very good behavior concerning memory consumption and computation time, as both depend only on the number of intervals  $I$  currently managed by the algorithm (see Figure 6(d)). In order to compute the position of a data item  $x$ , the strategy only needs to locate the interval containing  $f_B(x)$ , which can be done in  $O(\log I)$  using an appropriate tree structure. Furthermore, the algorithm strives to reduce the number of intervals created in each step in order to minimize memory consumption as much as possible. In practice, this yields an average access time of 5  $\mu s$  for a single data item and 13  $\mu s$  for 8 copies, while keeping a memory footprint similar to that of Redundant Share.

### C. Adaptivity

Adaptivity to changing environments is an important requirement for data distribution strategies and one of the main drawbacks of standard RAID approaches. Adding a single disk to a RAID system typically either requires the replacement of all data items in the system or splitting the RAID environment into multiple independent domains.

The theory behind randomized data distribution strategies claims that these strategies are able to compete with a best possible strategy in an adaptive setting. This means that the number of data movements to keep the properties of the strategy after a storage system has been inserted or deleted can be bounded against the best possible strategy. We assume in the following that a best possible algorithm just moves as much data from old disks to new disks, respectively from removed disks to remaining disks, as necessary to have the same usage on all storage systems. All bars in Figure 7 have been normalized to this definition of an optimal algorithm.

Furthermore, we distinguish between placements, where the ordering of the data items is relevant and where it is not. The first case occurs, e.g., for standard parity codes, where each data item has a different meaning (labeled "moved keeping order" in Figure 7). If a client accesses the third block of a parity set, then it is necessary to receive exactly that block. In contrast, the second case occurs for RAID 1 sets, where each copy has the same content and receiving any of this blocks is sufficient (labeled "moved changing order"). We will see in the following that not having to keep the order strongly simplifies the rebalancing process.

We start our tests in all cases with 128 storage systems and increase the number of storage systems by 1, 2, 3, 5, 7, 11, or 13 storage systems. The new storage systems have 1,5-times the capacity of the original system.

The original Consistent Hashing paper shows that the number of replacements is optimal for Consistent Hashing by showing that data is only moved from old disks to new disks in case of the insertion of a storage system or from a removed disk to old disks in the homogeneous setting [17]. Figure 7(a) shows a very different behavior, the number of data movements is sometimes more than 20-times higher than necessary. The reason is that we are placing  $400 \cdot \lceil \log n \rceil$  points

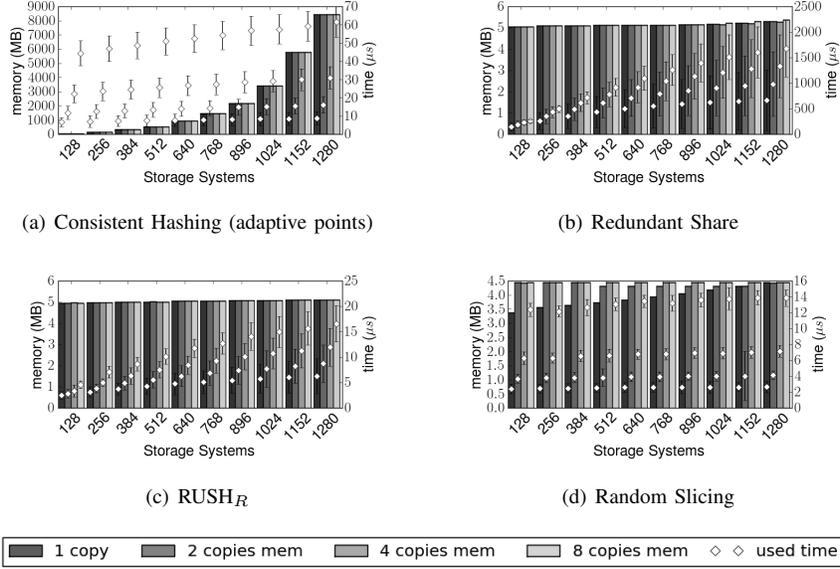


Fig. 6. Memory consumption and performance of the data distribution strategies in a heterogeneous setting.

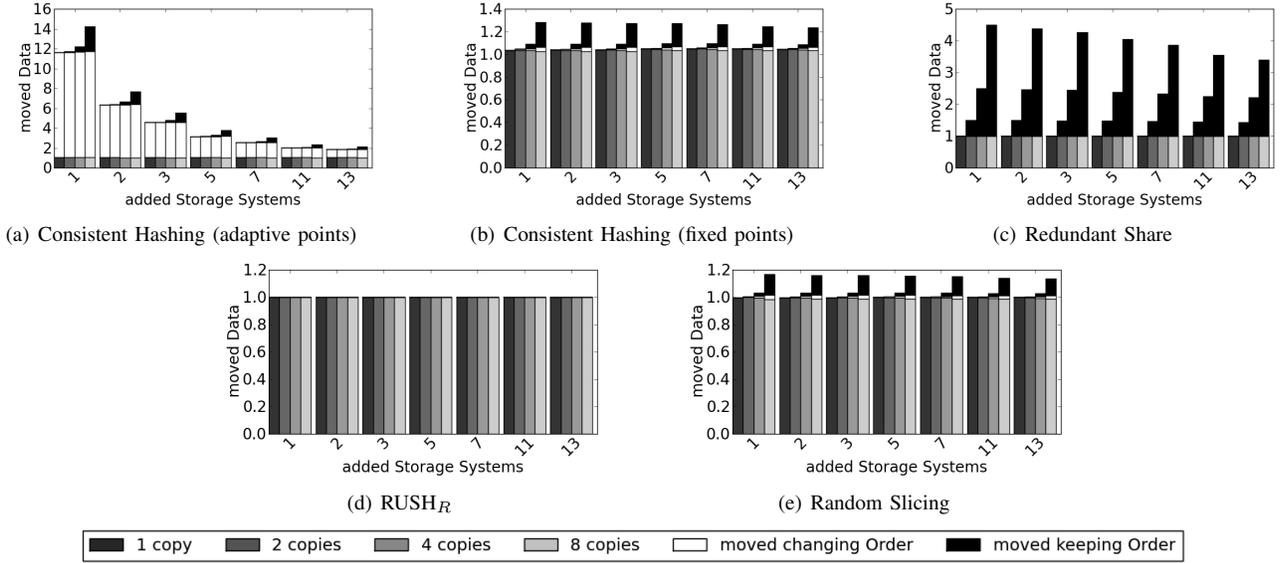


Fig. 7. Adaptivity of the data distribution strategies in a heterogeneous setting.

for each storage system and  $\lceil \log n \rceil$  increases from 7 to 8 when adding storage system number 129. This leads to a large number of data movements between already existing storage systems. Furthermore, the competitiveness strongly depends on whether the ordering of the different copies has to be maintained or not.

Figure 7(b) shows the adaptivity of Consistent Hashing in case that the number of points is fixed for each individual storage system and only depends on its own capacity. We use 2,400 points for the smallest storage system and use a proportional higher number of points for bigger storage systems. In this case the insertion of new storage systems only leads to data movements from old systems to the new

ones and not between old ones and therefore the adaptivity is very good in all cases. Figure 8 shows that the fairness in this case is still acceptable even in a heterogeneous setting.

The adaptivity of Redundant Share for adding new storage systems is nearly optimal, which is in line with the proofs presented in [5]. Nevertheless, Redundant Share is only able to achieve an optimal competitiveness if a new storage system is inserted that is at least as big as the previous ones. Otherwise it can happen that Redundant Share is only  $\log n$ -competitive (see Figure 7(c)).

Figure 7(d) shows that  $RUSH_R$  performs nearly optimal when storage nodes are added. Note, however, that we did not evaluate the effect on replica ordering because the current

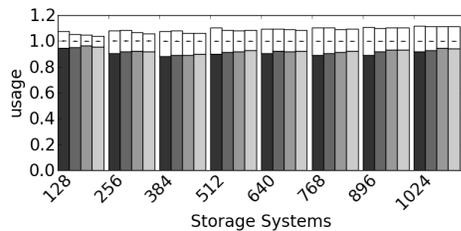


Fig. 8. Fairness of Consistent Hashing for fixed number of points in a heterogeneous setting.

implementation does not support replicas as distinct entities. Instead,  $RUSH_R$  distributes all replicas within one cluster.

Figure 7(e) shows that the adaptivity of Random Slicing is very good in all cases. This is explained because intervals for new storage systems are always created from fragments of old intervals, thus forcing data items to migrate only to new storage systems.

## VI. CONCLUSIONS

This paper shows that many randomized data distribution strategies are unable to scale to exascale environments, as either their memory consumption, their load deviation, or their processing overhead is too high. Nevertheless, they are able to easily adapt to changing environments, a property which cannot be delivered by table- or rule-based approaches.

The proposed Random Slicing strategy combines the advantages of all these approaches by keeping a small table and thereby reducing the amount of necessary random experiments. The presented evaluation and comparison with well-known strategies shows that Random Slicing is able to deliver the best fairness in all the cases studied and to scale up to exascale data centers.

## REFERENCES

- [1] S. Amarasinghe, D. Campbell, W. Carlson, A. Chien, W. Dally, E. Elnohazy, M. Hall, et al. Exascale software study: Software challenges in extreme scale systems. Technical report, sponsored by DARPA IPTO in the context of the ExaScale Computing Study, 2010.
- [2] A. Azagury, V. Dreizin, M. Factor, E. Henis, D. Naor, Y. Rinetzky, O. Rodeh, J. Satran, A. Tavory, and L. Yerushalmi. Towards an object store. In *Proceedings of the 20<sup>th</sup> IEEE Conference on Mass Storage Systems and Technologies (MSST)*, pages 165–176, 2003.
- [3] M. Blaum, J. Brady, J. Bruck, and J. Menon. EVENODD: an optimal scheme for tolerating double disk failures in RAID architectures. In *Proceedings of the 21<sup>st</sup> International Symposium on Computer Architecture (ISCA)*, pages 245–254, Apr. 1994.
- [4] A. Brinkmann and S. Effert. Redundant data placement strategies for cluster storage environments. In *Proceedings of the 12<sup>th</sup> International Conference On Principles Of Distributed Systems (OPODIS)*, Dec. 2008.
- [5] A. Brinkmann, S. Effert, F. Meyer auf der Heide, and C. Scheideler. Dynamic and Redundant Data Placement. In *Proceedings of the 27<sup>th</sup> IEEE International Conference on Distributed Computing Systems (ICDCS)*, Toronto, Canada, June 2007.
- [6] A. Brinkmann, M. Heidebuer, F. Meyer auf der Heide, U. Rückert, K. Salzwedel, and M. Vodisek. V: Drive - costs and benefits of an out-of-band storage virtualization system. In *Proceedings of the 21<sup>st</sup> IEEE Conference on Mass Storage Systems and Technologies (MSST)*, pages 153–157, 2004.
- [7] A. Brinkmann, K. Salzwedel, and C. Scheideler. Efficient, distributed data placement strategies for storage area networks. In *Proceedings of the 12<sup>th</sup> ACM Symposium on Parallel Algorithms and Architectures (SPAA)*, pages 119–128, 2000.

- [8] A. Brinkmann, K. Salzwedel, and C. Scheideler. Compact, adaptive placement schemes for non-uniform distribution requirements. In *Proceedings of the 14<sup>th</sup> ACM Symposium on Parallel Algorithms and Architectures (SPAA)*, pages 53–62, Winnipeg, Manitoba, Canada, 2002.
- [9] P. Corbett, B. English, A. Goel, T. Grcanac, S. Kleiman, J. Leong, and S. Sankar. Row-diagonal parity for double disk failure correction. In *Proceedings of the 3<sup>rd</sup> USENIX Conference on File and Storage Technologies (FAST)*, pages 1–14, San Francisco, CA, 2004.
- [10] T. Cortes and J. Labarta. Extending heterogeneity to RAID level 5. In *Proceedings of the USENIX Annual Technical Conference*, pages 119–132, Boston, Massachusetts, June 2001.
- [11] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Vosshall, and W. Vogels. Dynamo: Amazon’s highly available key-value store. *ACM SIGOPS Operating Systems Review*, 41(6):205–220, 2007.
- [12] A. Devulapalli, D. Dalessandro, and P. Wyckoff. Data Structure Consistency Using Atomic Operations in Storage Devices. In *Proceedings of the 5<sup>th</sup> International Workshop on Storage Network Architecture and Parallel I/Os (SNAPI)*, pages 65 – 73, Baltimore, USA, 2008.
- [13] J. Gonzalez and T. Cortes. Distributing Orthogonal Redundancy on Adaptive Disk Arrays. In *Proceedings of the International Conference on Grid computing, high-performance and Distributed Applications (GADA)*, Monterrey, Mexico, Nov. 2008.
- [14] R. J. Honicky and E. L. Miller. A fast algorithm for online placement and reorganization of replicated data. In *Proceedings of the 17<sup>th</sup> IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, Nice, France, Apr. 2003.
- [15] R. J. Honicky and E. L. Miller. Replication Under Scalable Hashing: A Family of Algorithms for Scalable Decentralized Data Distribution. In *Proceedings of the 18<sup>th</sup> IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2004.
- [16] N. L. Johnson and S. Kotz. *Urn Models and Their Applications*. John Wiley and Sons, New York, 1977.
- [17] D. Karger, E. Lehman, T. Leighton, M. Levine, D. Lewin, and R. Panigrahy. Consistent hashing and random trees: Distributed caching protocols for relieving hot spots on the World Wide Web. In *Proceedings of the 29<sup>th</sup> ACM Symposium on Theory of Computing (STOC)*, pages 654–663, El Paso, Texas, USA, 1997.
- [18] M. Mense and C. Scheideler. Spread: An adaptive scheme for redundant and fair storage in dynamic heterogeneous storage systems. In *Proceedings of the 19<sup>th</sup> ACM-SIAM Symposium on Discrete Algorithms (SODA)*, San Francisco, California, Jan. 2008.
- [19] M. Mitzenmacher. *The Power of Two Choices in Randomized Load Balancing*. PhD thesis, Computer Science Department, University of California at Berkeley, 1996.
- [20] D. A. Patterson, G. Gibson, and R. H. Katz. A Case for Redundant Arrays of Inexpensive Disks (RAID). In *Proceedings of the 1988 ACM Conference on Management of Data (SIGMOD)*, pages 109–116, June 1988.
- [21] P. Sanders. Reconciling simplicity and realism in parallel disk models. In *Proceedings of the 12<sup>th</sup> ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 67–76. SIAM, Philadelphia, PA, 2001.
- [22] C. Schindelhauer and G. Schomaker. Weighted distributed hash tables. In *Proceedings of the 17<sup>th</sup> ACM Symposium on Parallel Algorithms and Architectures (SPAA)*, pages 218–227, Las Vegas, Nevada, USA, July 2005.
- [23] I. Stoica, R. Morris, D. Liben-Nowell, D. Karger, M. Kaashoek, F. Dabek, and H. Balakrishnan. Chord: a scalable peer-to-peer lookup protocol for internet applications. *IEEE/ACM Transactions on Networking*, 11(1):17–32, 2003.
- [24] S. A. Weil, S. A. Brandt, E. L. Miller, D. D. E. Long, and C. Maltzahn. Ceph: A scalable, high-performance distributed file system. In *Proceedings of the 7<sup>th</sup> Symposium on Operating Systems Design and Implementation (OSDI)*, pages 307–320, Seattle, WA, USA, Nov. 2006.
- [25] S. A. Weil, S. A. Brandt, E. L. Miller, and C. Maltzahn. CRUSH: Controlled, Scalable And Decentralized Placement Of Replicated Data. In *Proceedings of the ACM/IEEE Conference on Supercomputing*, Tampa, FL, Nov. 2006.
- [26] W. Zheng and G. Zhang. FastScale: Accelerate raid scaling by minimizing data migration. In *Proceedings of the 9th USENIX Conference on File and Storage Technologies (FAST)*, Feb. 2011.