

On the Optimal Repair-Scaling Trade-off in Locally Repairable Codes

Si Wu, Zhirong Shen, and Patrick P. C. Lee

Department of Computer Science and Engineering, The Chinese University of Hong Kong

Abstract—How to improve the repair performance of erasure-coded storage is a critical issue for maintaining high reliability of modern large-scale storage systems. Locally repairable codes (LRC) are one popular family of repair-efficient erasure codes that mitigate the repair bandwidth and are deployed in practice. To adapt to the changing demands of access efficiency and fault tolerance, modern storage systems also conduct frequent scaling operations on erasure-coded data. In this paper, we analyze the optimal trade-off between the repair and scaling performance of LRC in clustered storage systems. Specifically, we design placement strategies that operate along the optimal repair-scaling trade-off curve subject to the fault tolerance constraints. We prototype and evaluate our placement strategies on a LAN testbed, and show that they outperform the conventional placement scheme in repair and scaling operations.

I. INTRODUCTION

As storage systems continue to scale, failures become commonplace [7]. To provide reliability guarantees for data storage even in the presence of failures, modern storage systems increasingly adopt *erasure coding* to maintain fault tolerance at low redundancy. Compared to traditional replication, erasure coding provably achieves higher degree of reliability at the same degree of redundancy [27], and has been widely deployed in enterprise storage systems [7], [11], [14], [21]. At a high level, erasure coding takes a set of original data blocks as input, and generates additional redundant blocks called *parity blocks*, such that all original data blocks can be reconstructed from a subset of available data and parity blocks.

To maintain data availability, storage systems need to perform frequent *repair* operations to recover any lost data from failures. However, erasure coding amplifies the network traffic and disk I/Os in repair operations [19]. Specifically, the repair of a single lost block incurs network transfers and disk I/Os of multiple available blocks for reconstruction. To mitigate the expensive repair cost introduced by erasure coding, many erasure code constructions have been proposed in the literature to improve repair efficiency. *Locally repairable codes (LRC)* [11], [14], [21], [24], in particular, are a new family of erasure codes that mitigate the network traffic and disk I/Os of repair with slight addition of storage redundancy. The main idea of LRC is to partition the original data blocks into several small-size local groups and generate a *local parity block* for each local group, such that the repair of any single failed block can be achieved within a short-size group. Due to its simplicity, ease of implementation, and the repair efficiency, LRC has been widely deployed in production [11], [14], [21].

In addition to repair, storage systems also perform frequent *scaling* operations to adapt to the requirements of fault tolerance and access efficiency. By scaling, we refer to the change of erasure coding parameters to balance the trade-off between access performance and storage efficiency [32]. Scaling is important for not only the expansion in storage capacity [34], but also adapting to the storage redundancy requirements with respect to the change of storage reliability [13].

In this paper, we analyze the interplay of both repair and scaling operations of LRC in *clustered storage systems*, which hierarchically organize nodes in multiple *clusters* such that the cross-cluster network bandwidth appears much more scarce than the inner-cluster bandwidth [2], [5], [26]. We show that the cross-cluster repair and scaling costs cannot be simultaneously minimized, and there exists a fundamental trade-off between the repair and scaling performance. To the best of our knowledge, this is the first work that unveils the inherent optimal repair-scaling trade-off in erasure-coded storage. To summarize, our contributions are as follows.

- We present a formal analysis on the optimal trade-off between the repair and scaling performance of LRC in clustered storage systems (Section III). We first derive the feasible data placement strategies subject to the single-cluster fault tolerance constraint, under which we explore different data placement strategies that operate along the optimal repair-scaling trade-off curve. Specifically, for a given (cross-cluster) scaling cost, we can find a data placement strategy that minimizes the (cross-cluster) repair cost.
- We implement the two extreme points of our placement strategies: one on minimizing the scaling cost, and another on minimizing the repair cost. We conduct experiments on a LAN testbed. Experimental results show that the placement with the minimum scaling cost reduces the scaling time of the baseline by up to 95.2%, while the placement with the minimum repair cost reduces the repair time of the baseline by up to 91.5% (Section IV).

II. BACKGROUND

A. Clustered Storage System Architecture

We consider a clustered storage system modeled as a two-level hierarchical architecture [5], as shown in Figure 1. The system partitions *storage nodes* (or *nodes* in short) into multiple *clusters*, such that the nodes within a cluster are connected via the same switch, while multiple clusters are interconnected via a *network core*. A cluster can refer to a rack [11], [21] or a data

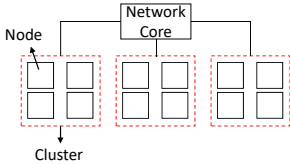


Fig. 1. Example of a clustered storage system architecture.

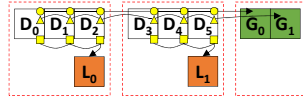


Fig. 2. LRC with $(k, l, g, c) = (6, 2, 2, 3)$, where each dotted-line represents a cluster.

center [3], [17]. The system organizes data in a collection of fixed-size *blocks*, which form the basic read/write units. Here, we assume that the cross-cluster bandwidth in clustered storage systems is much more scarce than the inner-cluster bandwidth [2], [5], [26]. Our goal is to minimize the cross-cluster network traffic in both repair and scaling operations.

B. Locally Repairable Codes

Basics of LRC. We present the basics of *locally repairable codes (LRC)*. We configure an LRC construction with four parameters (k, l, g, c) , meaning that there are k data blocks (denoted by D_0, D_1, \dots, D_{k-1}), l local parity blocks (denoted by L_0, L_1, \dots, L_{l-1}), and g global parity blocks (denoted by G_0, G_1, \dots, G_{g-1}), such that all $k + l + g$ data/parity blocks are stored in $k + l + g$ nodes located in c clusters. We call the set of $k + l + g$ data/parity blocks that are encoded together to be a *stripe*, and a large-scale storage system typically contains multiple stripes that are independently encoded. In the paper, our analysis focuses on a single stripe.

By storing each LRC stripe of blocks in multiple nodes and clusters, a storage system achieves both node-level and cluster-level fault tolerance. Since cluster failures are much less common than node failures in practice [17], we enforce the data placement of each LRC stripe such that the storage system provides only *single-cluster fault tolerance* (i.e., the data remains available under a single-cluster failure), while still providing multi-node fault tolerance. By storing multiple blocks in a single cluster, we can significantly reduce cross-cluster network traffic during repair [9], [10], [16], [18], [22].

There are various LRC constructions in the literature [11], [21], [24]. In this paper, we focus on the LRC construction based on Azure’s Local Reconstruction Codes [11]. Specifically, LRC divides the k data blocks into l equal-size groups, assuming that k is divisible by l . It computes an XOR sum based on each group of $\frac{k}{l}$ data blocks to form a local parity block. It also computes the g global parity blocks from all k data blocks. Let b be the number of data blocks that encode into a local parity block, i.e., $b = \frac{k}{l}$. We call the collection of b data blocks and their corresponding local parity block to be a *local group*. Figure 2 gives an example of LRC for $(k, l, g, c) = (6, 2, 2, 3)$.

Repair of LRC. LRC is designed to mitigate the network traffic and disk I/Os for a single-block repair (which is much more common than a multi-block repair in practice [11], [19]) by limiting the repair within a local group. Given that transient failures account for the majority of failure events [7] and the temporarily unavailable data blocks are served by

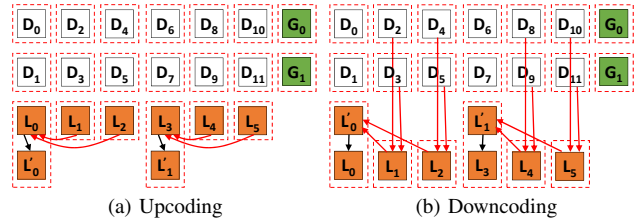


Fig. 3. Upcoding and downcoding examples for fast LRC $(12, 6, 2, 20)$ and compact LRC $(12, 2, 2, 16)$, in which each block is stored in a distinct cluster.

degraded reads, in this paper, we focus on the repair of a single unavailable data block in LRC.

We quantify the (single-data-block) repair performance of LRC in clustered storage systems as follows. To repair any unavailable data block, LRC retrieves other available blocks of the same local group. Let $C(D_i)$ denote the amount of cross-cluster network traffic being transferred for repairing a data block D_i . Then we define the *repair cost* of LRC (a.k.a. the average repair cost in [14]) as $\frac{1}{k} \sum_{i=0}^{k-1} C(D_i)$.

If each local group is put in the same cluster (e.g., in Figure 2), the repair cost is zero as there is no cross-cluster network transfer. If LRC adopts the flat data placement by storing each of the $k + l + g$ blocks of a stripe in a distinct cluster [11], [21], [32], then the repair cost is equal to b .

Scaling of LRC. We define *scaling* as the change of coding parameters in erasure coding. In this paper, we consider the change of two sets of LRC parameters in response to workload changes as described in [32]. Specifically, we consider two LRC constructions: (i) *fast LRC*, which stores more local parity blocks for higher repair performance (e.g., for hot data), and (ii) *compact LRC*, which stores fewer local parity blocks for higher storage efficiency (e.g., for cold data). Both fast and compact LRCs store the same numbers of data blocks and global parity blocks. We perform scaling operations to switch between the fast LRC and the compact LRC to balance between access performance and storage efficiency. We call the conversion *upcoding* when scaling from the fast LRC to the compact LRC, and *downcoding* when scaling from the compact LRC to the fast LRC. Accordingly, we define the *upcoding cost* and the *downcoding cost* as the amounts of cross-cluster network traffic transferred for upcoding and downcoding, respectively.

Figure 3 depicts the upcoding and downcoding for the fast LRC with $(k, l, g, c) = (12, 6, 2, 20)$ and the compact LRC with $(k, l, g, c) = (12, 2, 2, 16)$. Both codes are deployed under the flat data placement (i.e., each block is stored in a distinct cluster). The upcoding operation (from $(12, 6, 2, 20)$ to $(12, 2, 2, 16)$) sends L_1 and L_2 across clusters to L_0 to compute a new local parity block L'_0 , and sends L_4 and L_5 across clusters to L_3 to compute a new local parity block L'_1 . Thus, the upcoding cost is 4. On the other hand, the downcoding operation (from $(12, 2, 2, 16)$ to $(12, 6, 2, 20)$) sends D_2 and D_3 to another cluster to compute L_1 , and sends D_4 and D_5 to another cluster to compute L_2 . It then sends L_1 and L_2 to the cluster that holds L'_0 to compute L_0 . The same is for L_3, L_4 , and L_5 . In total, the downcoding cost is 12.

III. TRADE-OFF BETWEEN REPAIR AND SCALING

We now show that there exists a fundamental trade-off between the repair and scaling (i.e., upcoding/downcoding) performance in clustered storage systems.

A. Fault Tolerance of LRC

We first analyze the feasible data placement strategies subject to the single-cluster fault tolerance constraint (Section II-B).

Lemma 1. *A (k, l, g) LRC can tolerate any $g + i$ block failures that span i local groups, where $1 \leq i \leq l$. However, the fault tolerance will fail if there exist i local groups with more than $g + i$ block failures.*

Proof. Consider a set of up to $g + i$ failed blocks that span i local groups ($1 \leq i \leq l$). Suppose that the failed blocks comprise d data blocks, x local parity blocks ($0 \leq x \leq i$), and y global parity blocks ($0 \leq y \leq g$), such that $d + x + y \leq g + i$. In other words, x local groups have failed local parity blocks, while $i - x$ local groups have failed data blocks only and their local parity blocks are available for repair. For each of the $i - x$ local groups with available local parity blocks, we can swap the available local parity block with one failed data block, and mark the data block as available and the local parity block as failed. By doing so, we can transform the set of failed blocks into $d - (i - x)$ failed data blocks, which can now be decoded by the $g - y$ surviving global parity blocks as $d - (i - x) \leq g - y$ [11]. After decoding the failed data blocks, we can recover all failed local/global parity blocks. The fault tolerance is maintained.

Based on the above analysis, we can also prove that a (k, l, g) LRC cannot tolerate any more than $g + i$ block failures if they span i local groups ($1 \leq i \leq l$). By swapping the available local parity blocks and the failed data blocks, we have $d - (i - x)$ failed data blocks that need to be decoded from $g - y$ available global parity blocks. If there are more than $g + i$ block failures, we have $d - (i - x) > g - y$, implying that there are more failed data blocks than the available global parity blocks. Thus, the data blocks cannot be repaired. \square

By Lemma 1, we can deduce that if we provide single-cluster fault tolerance for a (k, l, g) LRC, we have to place no more than $g + i$ blocks that span i local groups in a cluster; otherwise, a cluster failure will cause data loss.

B. Motivating Examples

We show via motivating examples that there is no data placement that can simultaneously minimize both the repair and scaling costs for LRC. Here, we focus on comparing the repair of the fast LRC and the upcoding from the fast LRC (for hot data) to the compact LRC (for cold data), since the transition from hot data to cold data is more common in real-world storage workloads [17].

Figure 4 presents two data placements for a fast LRC with $(12, 6, 2, 7)$, which is to be upcoded to a compact LRC with $(12, 2, 2, 7)$. As l reduces from 6 to 2, we generate a local parity block of the compact LRC based on every three local parity blocks of the fast LRC.

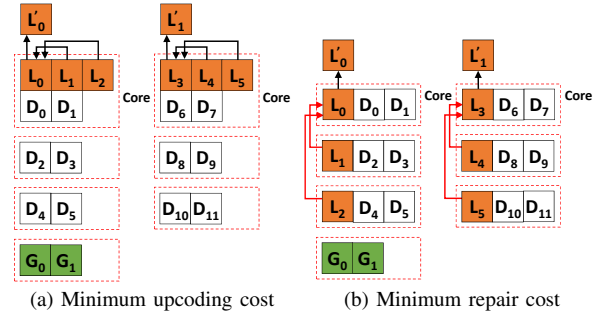


Fig. 4. Two data placements for the fast LRC with $(k, l, g, c) = (12, 6, 2, 7)$ being upcoded to the compact LRC with $(k, l, g, c) = (12, 2, 2, 7)$.

In Figure 4(a), we place the local parity blocks in two clusters. Since there is no cross-cluster transfer in upcoding, the upcoding cost is zero, which is the minimum. On the other hand, the cost for repairing each of D_0 to D_5 is 0, 0, 1, 1, 1, 1, respectively, while the same is for D_6 to D_{11} . Thus, the repair cost is 0.67.

In Figure 4(b), we place each local group (the local parity block and its corresponding encoding data blocks) in one cluster. The repair of each data block can now be done in each cluster locally, so the repair cost is zero, which is the minimum. However, the upcoding needs to transfer four blocks across clusters. Thus, the upcoding cost is four.

Figure 4 shows that the repair and upcoding costs cannot be simultaneously minimized for any possible data placement.

C. Trade-off Analysis

Preliminaries. We now formally analyze the trade-off between repair and upcoding (we address the case of downcoding in Section III-D). We first present the definitions and notations. Since the access performance of LRC is mainly related to the number of local parity blocks, we consider the scaling operation that varies the number of local parity blocks as in [32]. To simplify our analysis, we consider the case where the number of local parity blocks for the fast LRC (denoted by l) is divisible by the number of local parity blocks for the compact LRC (denoted by l'), such that every local parity block of the compact LRC can be updated from $\frac{l}{l'}$ local parity blocks of the fast LRC. Let δ be the *scaling factor*, defined as $\delta = \frac{l}{l'}$.

For a local parity block L_i , we call the set of data blocks that generates L_i a *local data set*, denoted by $\mathcal{E}_i = \{D_{i \times b}, \dots, D_{(i+1) \times b - 1}\}$, where b is the number of data blocks that are encoded to L_i (defined in Section II-B). For example, in Figure 4, the local data sets for L_0 and L_1 are $\mathcal{E}_0 = \{D_0, D_1\}$ and $\mathcal{E}_1 = \{D_2, D_3\}$, respectively. Note that \mathcal{E}_i and L_i together form a local group. During upcoding, we convert every δ local groups of the fast LRC into one local group of the compact LRC, and we call these δ local groups an *upcoding unit*. There are a total of l' upcoding units, and the i -th one is composed of the blocks $\mathcal{E}_{i \times \delta}, L_{i \times \delta}, \dots, \mathcal{E}_{(i+1) \times \delta - 1}, L_{(i+1) \times \delta - 1}$. For example, in Figure 4, we have $\delta = 3$, and there are two upcoding units: (i) $\{\mathcal{E}_0, L_0, \mathcal{E}_1, L_1, \mathcal{E}_2, L_2\}$, and (ii) $\{\mathcal{E}_3, L_3, \mathcal{E}_4, L_4, \mathcal{E}_5, L_5\}$.

For each upcoding unit, we define a *core cluster* as the cluster that stores the δ local parity blocks and aggregates them into one local parity block of the compact LRC. For example, in Figure 4(a), a core cluster stores L_0, L_1 , and L_2 and encodes them into L'_0 , while another core cluster stores L_3, L_4 , and L_5 and encodes them into L'_1 . Suppose that $b \leq g$, such that a local data set can be entirely stored in one cluster without breaking single-cluster fault tolerance. Let θ be the maximum number of local data sets that can be collocated with their corresponding local parity blocks in one cluster. By Lemma 1, the number of data and local parity blocks (i.e., $\theta \times b + \theta$), which span θ local groups, cannot exceed $g + \theta$. Thus, we can calculate θ as $\theta = \lfloor \frac{g}{b} \rfloor$.

For example, in Figure 4, every $\theta = 1$ local data set can be collocated with its local parity block in one cluster. Our analysis focuses on $b \leq g$, while we will later show that the analysis for $b > g$ is similar.

Roadmap. Our trade-off analysis between repair and upcoding is organized as follows. (i) We first design the placement of the local parity blocks to achieve the globally minimum upcoding cost. (ii) Given the condition that the upcoding cost is minimized, we vary the locations of the data blocks to minimize the repair cost; note that the repair cost is not necessarily globally minimum (Section III-B). (iii) We gradually relocate the local parity blocks, so as to trade the increased upcoding cost for the decreased repair cost and finally achieve the globally minimum repair cost. The number of clusters (i.e., c) is determined by the placement policy. Since the global parity blocks do not participate in repair and scaling operations, we simply place them in a dedicated cluster and omit their discussion in our analysis.

Guiding example. For better understanding of our analysis, we use a guiding example that upcodes from the fast LRC with $(12, 6, 2, 7)$ to the compact LRC with $(12, 2, 2, 7)$.

- *Minimizing upcoding cost:* First, from Figure 5(a), we place L_0, L_1 , and L_2 in a core cluster, and place L_3, L_4 , and L_5 in another core cluster. We perform upcoding in each core cluster by aggregating the three local parity blocks without any cross-cluster transfer. The upcoding cost is zero.
- *Minimizing repair cost under minimum upcoding cost:* Next, in Figure 5(b), we place \mathcal{E}_0 in one core cluster, and place \mathcal{E}_1 and \mathcal{E}_2 in two different clusters. We place $\mathcal{E}_3, \mathcal{E}_4$, and \mathcal{E}_5 in the same way. We can show that the cost for repairing each data block in \mathcal{E}_0 and \mathcal{E}_3 is zero, while the cost for repairing each data block in $\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_4$, and \mathcal{E}_5 is one. The repair cost of the fast LRC is $\frac{4}{6} = 0.67$. We will show that this repair cost is minimized under the minimum upcoding cost.
- *Trading increased upcoding cost for decreased repair cost:* In this example, we see that the main reason that causes non-zero repair cost is the separate placement of the local parity block and its corresponding local data set. For example, in Figure 5(b), L_1 is stored in the core cluster, while \mathcal{E}_1 is stored in a different cluster; the same holds for L_2, L_4 , and L_5 . Thus, we can gradually move one local parity block from the core cluster to the cluster where its local data set resides.

Algorithm 1 Trade-off placement

Input: Integer x ($0 \leq x \leq (\delta - \theta) \times l'$), which is a multiple of θ
Output: A placement for the fast LRC

```

1: for the  $i$ -th ( $0 \leq i \leq l' - 1$ ) upcoding unit do
2:   // Minimizing upcoding cost
3:   Select a new core cluster
4:   Put  $L_{i \times \delta}, \dots, L_{(i+1) \times \delta - 1}$  in the core cluster
5:   // Further minimizing repair cost
6:   for  $j = 0$  to  $\theta - 1$  do
7:     Put  $\mathcal{E}_{i \times \delta + j}$  into the core cluster
8:   end for
9:   for  $j = \theta$  to  $\delta - 1$  do
10:    if  $j \bmod \theta = 0$  then
11:      Select a new cluster
12:    end if
13:    Put  $\mathcal{E}_{i \times \delta + j}$  into the new cluster
14:  end for
15: end for
16: // Trading upcoding cost for repair cost
17: Move  $x$  local parity blocks from the core clusters to the clusters
    where their corresponding local data sets reside

```

In Figure 5(c), we move L_1 to the cluster that holds \mathcal{E}_1 . By doing so, we reduce the repair cost by $\frac{1}{6} = 0.17$, while the upcoding cost increases by one (the core cluster must now retrieve the relocated L_1 across cluster for upcoding). We will also show that the repair cost in Figure 5(c) (i.e., $\frac{3}{6} = 0.5$) is minimized subject to the upcoding cost of one.

- *Minimizing repair cost:* Finally, in Figure 5(d), we move all of L_1, L_2, L_4 , and L_5 to the clusters where their corresponding local data sets reside. As the repair can now be executed without any cross-cluster traffic, the repair cost is zero.

Algorithm design. We now present a data placement algorithm (Algorithm 1) that is guaranteed to operate on the optimal trade-off curve from the minimum upcoding cost to the minimum repair cost. The input of Algorithm 1 is a parameter x that decides the operation point in the optimal trade-off curve, while the output is a placement for the fast LRC.

(i) *Placing local parity blocks to minimize upcoding cost to zero.* We first gather the δ local parity blocks into the selected core cluster in each upcoding unit (lines 3-4), such that we can complete upcoding within each core cluster. By doing this, the upcoding cost is zero. For example, in Figure 5(a), we store every $\delta = 3$ local parity blocks in each core cluster.

(ii) *Placing data blocks to further minimize repair cost ($x = 0$).* We next determine the locations of the data blocks to minimize the repair cost given the condition that the upcoding cost has been minimized. We first focus on an upcoding unit. There are δ local parity blocks that are collocated in the core cluster, and we have to decide how to place the δ local data sets. If we put more local data sets also in the core cluster, then more data blocks can be repaired in the core cluster locally and the repair cost can be minimized. Note that the blocks in the core cluster span δ local groups, so the sum of the number of local parity blocks (i.e., δ) and the number of data blocks cannot exceed $g + \delta$ to promise single-cluster fault tolerance. As a result, we can put $\theta = \lfloor \frac{g}{b} \rfloor$ local data sets in the core

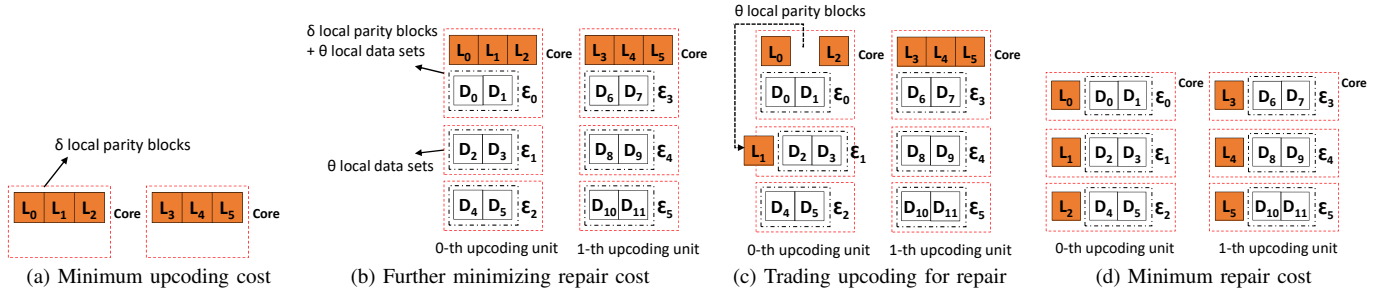


Fig. 5. Illustration of the different steps of our placement policy for fast LRC with $(k, l, g, c) = (12, 6, 2, 7)$ scaled to compact LRC with $(k, l, g, c) = (12, 2, 2, 7)$.

cluster (lines 6-8), and the cost for repairing each data block in these θ local data sets is zero. For example, in Figure 5(b), we put one local data set (i.e., $\theta = 1$) into a core cluster.

For the remaining $\delta - \theta$ local data sets, We make sure that each local data set is entirely stored in one cluster (not the core cluster), such that the cost for repairing each data block in these $\delta - \theta$ local data sets is one, as the corresponding local parity block is in the core cluster. In this manner, the repair cost is minimized. We further collocate every θ local data sets into one different cluster (lines 9-14). Each such cluster will have $\theta \times b = \lfloor \frac{g}{b} \rfloor \times b \leq g$ data blocks, hence complying with single-cluster fault tolerance. For example, in Figure 5(b), for the remaining two local data sets, we put each into a different cluster. The data blocks of other upcoding units are placed in the same way. We call this placement *Opt-S*, where the upcoding cost is zero and the repair cost is $\frac{\delta - \theta}{\delta} = 1 - \frac{\theta}{\delta}$. Note that *Opt-S* is derived by inputting $x = 0$ in Algorithm 1.

The values of δ and θ significantly influence the upcoding and repair costs, which we discuss as follows.

- If $\delta \leq \theta$, then we can directly put all δ local data sets of an upcoding unit in the core cluster. Under this placement, the repair and upcoding operations can be directly performed within the core cluster, and therefore the repair and upcoding costs are zero.
- In the case where $\delta > \theta$, for the remaining $\delta - \theta$ local data sets, we collocate every θ ones into a different cluster. If we collocate more than θ local data sets in a different cluster, then when we move the corresponding local parity blocks into this cluster, it will violate single-cluster fault tolerance. If we collocate less than θ , say s ($s < \theta$) local data sets in a different cluster, then after we relocate the s corresponding local parity blocks into this cluster, the repair cost reduces by $\frac{s}{\theta}$ while the upcoding cost increases by one. However, in our design, after we relocate the θ corresponding local parity blocks into this cluster, the repair cost reduces by $\frac{\theta}{\theta}$ while the upcoding cost increases by one. That is to say, the reduction of the repair cost is the most in our design while the upcoding cost increases by one.
- We assume that $\delta - \theta$ is divisible by θ , such that every θ local data sets can be stored in the selected cluster.

(iii) *Relocating local parity blocks to trade upcoding cost for repair cost* ($0 \leq x \leq (\delta - \theta) \times l'$). In *Opt-S*, for each upcoding unit, there are $\delta - \theta$ local parity blocks lying in the core cluster,

while their local data sets are located in different clusters. If we move such local parity blocks to where their local data sets reside, then the repair cost can be further reduced.

Since we collocate every θ local data sets into one different cluster, we can move θ corresponding local parity blocks to this cluster. By doing this, we reduce the repair cost by $\frac{\theta}{\theta}$ as we move θ local parity blocks to be collocated with their local data sets, such that the cost for repairing any data block in these θ local data sets reduces from one to zero. During upcoding, we first apply partial encoding of the θ local parity blocks to calculate an XOR sum in this cluster, and then send the XOR sum to the core cluster. Thus, the upcoding cost increases by one. We move (θ local parity blocks) for one different cluster in a step, and in an *upcoding unit by upcoding unit* basis, to transform *Opt-S* into a placement that trades increased upcoding cost for decreased repair cost. Since there are $\delta - \theta$ local parity blocks (per upcoding unit) that can be moved and l' upcoding units, we can then move x local parity blocks for $\frac{x}{\theta}$ different clusters, where x is a multiple of θ and $0 \leq x \leq (\delta - \theta) \times l'$ (line 17). The repair cost reduces by $\frac{x}{\theta}$ and the upcoding cost increases by $\frac{x}{\theta}$ compared to those costs of *Opt-S*. Thus, the cost values of the placement derived in a transformation step are shown as follows.

$$\begin{aligned} \text{upcoding cost} &= \frac{x}{\theta} \\ \text{repair cost} &= 1 - \frac{\theta}{\delta} - \frac{x}{\theta}. \end{aligned} \quad (1)$$

For example, in Figure 5(c), we move L_1 to the cluster that holds \mathcal{E}_1 , and we reduce the repair cost by $\frac{1}{6} = 0.17$ while the upcoding cost increases by one.

Note that each transformation step will not break the fault tolerance guarantee. Specifically, suppose that a core cluster has moved out t (t is a multiple of θ and $0 \leq t \leq \delta - \theta$) local parity blocks, it then remains $\delta - t$ local parity blocks and θ local data sets. The number of blocks is $(\theta \times b + \delta - t) \leq (g + \delta - t)$, and the blocks span $\delta - t$ local groups. For a different cluster that accommodates the relocated local parity blocks, it now has θ local parity blocks collocated with their corresponding local data sets. According to Lemma 1, both clusters can guarantee single-cluster fault tolerance.

The movements (line 17) guarantee that the repair cost of the placement derived in a transformation step (i.e., $1 - \frac{\theta}{\delta} - \frac{x}{\theta}$) is minimized given the upcoding cost (i.e., $\frac{x}{\theta}$), which can be readily deduced by the following Theorem.

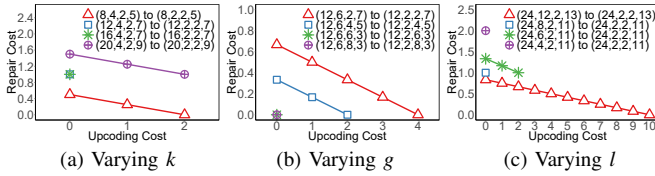


Fig. 6. Trade-off curve between the upcoding cost and the repair cost.

Theorem 1. For any placement subject to single-cluster fault tolerance, if the upcoding cost is u , then the lower bound of the repair cost is $1 - \frac{\theta}{\delta} - \frac{u \times \theta}{l}$.

The proof is elaborated in the Appendix. Since the repair cost of the placement derived in each transformation step touches the lower bound, it is minimized under the upcoding cost.

(iv) *Minimizing repair cost to zero* ($x = (\delta - \theta) \times l'$). In the end (i.e., inputting $x = (\delta - \theta) \times l'$ in Algorithm 1), every θ local parity blocks lie together with their local data sets in one cluster. We call this placement *Opt-R*, where all repair operations can be done within each cluster locally, so the repair cost is zero. For example, in Figure 5(d), every local parity block lies together with its local data set in one cluster.

Extension to $b > g$. We now demonstrate that the analysis for $b > g$ is similar. If $b \geq g + 1$ and we further assume that $b \bmod (g + 1) \neq 0$. Since at most $g + 1$ data blocks of a local data set can be put into one cluster without breaking the fault tolerance guarantee, a local data set can be put into $\lfloor \frac{b}{g+1} \rfloor + 1$ clusters, where $\lfloor \frac{b}{g+1} \rfloor$ ones hold $g + 1$ data blocks each, and the remaining one holds $(m = b \bmod (g + 1)) \leq g$ data blocks. If we focus on the remaining m data blocks, then the above analysis directly applies. The only change is that the repair cost has to be added by a difference of $\lfloor \frac{b}{g+1} \rfloor$ as we have to retrieve an XOR sum of all blocks in each of the other $\lfloor \frac{b}{g+1} \rfloor$ clusters to repair any data block in a local data set.

If b is divisible by $(g + 1)$, then a local data set spans $\frac{b}{g+1}$ clusters with $g + 1$ data blocks each. To achieve the minimum upcoding cost, we gather every δ local parity blocks in each core cluster, and then we can put at most another g data blocks in a core cluster (Lemma 1). We cannot move data blocks to reduce the number of clusters (except the core cluster) a local data set spans. As a result, the cost for repairing any data block cannot be further reduced. Thus, in the placement with the minimum upcoding cost, the repair cost is also minimized.

Trade-off exemplification. We plot the upcoding cost and the repair cost of each placement and obtain a trade-off curve between the upcoding cost and the repair cost. We give some examples in Figure 6 to show the trade-off curve. Several findings are stated as follows.

The most important trend is that as the upcoding cost increases, the repair cost decreases.

As shown in Figures 6(a) and 6(c), increasing k or decreasing l makes the overall repair cost increased. This is because b increases, which further results in larger $\lfloor \frac{b}{g+1} \rfloor$. For example, in Figure 6(a), in the case where $(k, l, g, c) = (8, 4, 2, 5)$, $\lfloor \frac{b}{g+1} \rfloor = 0$ (i.e., a local data set is entirely stored in one cluster), while

in the case where $(k, l, g, c) = (20, 4, 2, 9)$, $\lfloor \frac{b}{g+1} \rfloor = 1$ (i.e., a local data set has to span two clusters).

Figure 6(b) tells that, as g increases, there are less trade-off points, and the overall repair cost decreases. The reason is that larger g means that more local data sets can be put into the core cluster, such that the cost for repairing any data block therein is zero. For example, in the case where $(k, l, g, c) = (12, 6, 2, 7)$, we can put one local data set into the core cluster, while in the case where $(k, l, g, c) = (12, 6, 4, 5)$, we can put two local data sets into the core cluster.

In Figure 6(a), the case where $(k, l, g, c) = (12, 4, 2, 7)$ only exhibits one point as b is divisible by $(g + 1)$, and the case where $(k, l, g, c) = (16, 4, 2, 7)$ also has one point because $\delta \leq \lfloor \frac{g}{m} \rfloor$ where $m = b \bmod (g + 1)$. The cases with one point in Figure 6(b) are due to that $\delta \leq (\theta = \lfloor \frac{g}{b} \rfloor)$, and those in Figure 6(c) are due to that b is divisible by $(g + 1)$.

D. Downcoding Procedures

We now address the downcoding procedure to restore the original layout of the fast LRC. We mainly consider the downcoding processes for Opt-S and Opt-R. We can see that the block layout for each upcoding unit is the same in Opt-S and Opt-R, and we can readily derive that the downcoding operation for each upcoding unit is the same. Therefore, we conduct our analysis within a single upcoding unit.

Downcoding for Opt-S. In Opt-S, after upcoding, the δ local parity blocks in the core cluster are updated into a local parity block of the compact LRC. Downcoding is to rebuild the δ local parity blocks in the core cluster. Since there are θ local data sets (i.e., $\mathcal{E}_0, \dots, \mathcal{E}_{\theta-1}$) in the core cluster, we can first recalculate θ local parity blocks (i.e., $L_0, \dots, L_{\theta-1}$) within the core cluster. Note that there is a local parity block of the compact LRC (i.e., L'_0) in the core cluster that can be utilized, we then have to retrieve at least $\delta - \theta - 1$ local parity blocks from other clusters. Thus, for each of $\mathcal{E}_\theta, \dots, \mathcal{E}_{\delta-2}$, we locate the cluster that holds it, calculate the corresponding local parity block, and send the local parity block to the core cluster. Finally, we use $L_0, \dots, L_{\delta-2}$ and L'_0 to calculate $L_{\delta-1}$.

The downcoding cost is readily deduced as $(\delta - \theta - 1) \times l'$, if $b \leq g$. If considering general parameters, then when calculating each of the $\delta - 1$ local parity blocks, we should access partial encoded blocks from another $\lfloor \frac{b}{g+1} \rfloor$ clusters. Therefore,

$$\text{downcoding cost} = ((\delta - \lfloor \frac{g}{m} \rfloor - 1) + \lfloor \frac{b}{g+1} \rfloor \times (\delta - 1)) \times l', \quad (2)$$

where $m = b \bmod (g + 1)$.

For example, in Figure 4(a), during downcoding, we can recalculate L_0 based on \mathcal{E}_0 in the core cluster. We then calculate L_1 based on \mathcal{E}_1 in a different cluster, and send L_1 to the core cluster. Finally, we use L_0, L_1 and L'_0 to calculate L_2 . The same is for L_3, L_4 and L_5 , and thus the downcoding cost is two.

Downcoding for Opt-R. In Opt-R, downcoding is to restore the layout that every θ local parity blocks and their local data sets are collocated in one cluster. Thus, we can regenerate each local parity block within each cluster locally, implying that the downcoding cost is zero. For general parameters, when

recalculating each local parity block, we should access partial encoded blocks from another $\lfloor \frac{b}{g+1} \rfloor$ clusters. Thus,

$$\text{downcoding cost} = \lfloor \frac{b}{g+1} \rfloor \times \delta \times l' = \lfloor \frac{b}{g+1} \rfloor \times l. \quad (3)$$

For example, in Figure 4(b), we can regenerate L_0-L_5 within each cluster locally, so the downcoding cost is zero.

IV. EVALUATION

We evaluate via numerical analysis and testbed experiments the two optimal placements (Section III-C): (i) Opt-S, the placement with the minimum upcoding cost, and (ii) Opt-R, the placement with the minimum repair cost. We compare them with the *flat placement* (Flat) that stores each block of a stripe in a distinct cluster [11], [21].

A. Numerical Analysis

Our numerical analysis considers four metrics: the repair cost of the fast LRC, the repair cost of the compact LRC, the upcoding cost, and the downcoding cost.

Flat. For Flat, the repair costs of the fast LRC and the compact LRC are $\frac{k}{l} = b$ and $\frac{k}{l'} = b'$, respectively (Section II-B).

During upcoding, we send every $\delta - 1$ local parity blocks to the cluster that holds one remaining local parity block to encode into a local parity block of the compact LRC (per upcoding unit). Thus, the upcoding cost is calculated as

$$\text{upcoding cost} = (\delta - 1) \times l' = l - l'.$$

During downcoding, each of the first $\delta - 1$ local parity blocks is recalculated by sending its encoding data blocks across clusters to another cluster. The remaining one local parity block can be derived in two ways: (i) the same as the first $\delta - 1$ local parity blocks, and (ii) sending the $\delta - 1$ available local parity blocks to the cluster that holds the local parity block of the compact LRC. Hence,

$$\begin{aligned} \text{downcoding cost} &= \min\{\delta \times b \times l', ((\delta - 1) \times b + \delta - 1) \times l'\} \\ &= \min\{k, k + l - l' \times (b + 1)\}. \end{aligned}$$

Opt-S and Opt-R. The repair cost of the fast LRC of Opt-S and Opt-R is calculated using Equation (1), and the repair cost of the compact LRC of Opt-S and Opt-R is easily derived based on the layout after upcoding (Figure 4).

The upcoding cost of Opt-S and Opt-R is calculated using Equation (1). The downcoding cost of Opt-S is calculated using Equation (2), and that of Opt-R is calculated using Equation (3).

Scaling operations. We consider eight sets of scaling operations from (k, l, g, c) to (k', l', g, c) , denoted by p_0 to p_7 :

$$\begin{aligned} p_0: (4, 2, 2, 3) &\leftrightarrow (4, 1, 2, 3) & p_1: (8, 4, 2, 5) &\leftrightarrow (8, 2, 2, 5) \\ p_2: (12, 6, 2, 7) &\leftrightarrow (12, 2, 2, 7) & p_3: (12, 6, 4, 5) &\leftrightarrow (12, 2, 4, 5) \\ p_4: (24, 12, 4, 7) &\leftrightarrow (24, 2, 4, 7) & p_5: (24, 6, 2, 11) &\leftrightarrow (24, 2, 2, 11) \\ p_6: (32, 16, 4, 9) &\leftrightarrow (32, 4, 4, 9) & p_7: (32, 8, 2, 13) &\leftrightarrow (32, 2, 2, 13) \end{aligned}$$

Results. Figure 7 plots the costs for the three placements. We summarize the following observations.

- From Figures 7(a) and 7(c), Opt-S always has the minimum upcoding cost (zero), while Opt-R is designed with the minimum repair cost.

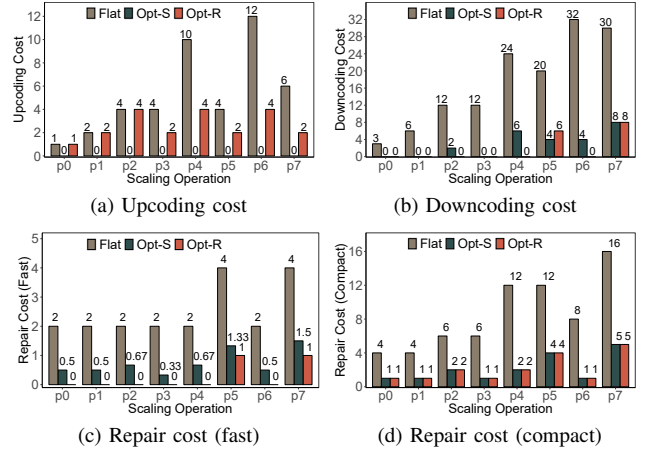


Fig. 7. Numerical results of the scaling and repair costs.

- Opt-S and Opt-R outperform Flat in terms of the scaling and repair costs stably.
- For the cases where $b = g$ (e.g., p_0-p_2), we can only put $\theta = 1$ local parity block and its local data set in one cluster, and so we cannot utilize partial encoding of the local parity blocks in each cluster for upcoding. Thus, the upcoding cost of Opt-R equals that of Flat. For other cases, we put more than one local parity block in one cluster, and so partial encoding of the local parity blocks brings benefit to the upcoding cost. For example, for scaling operation p_4 , the upcoding cost of Flat is 10 while that of Opt-R is four.
- In Figure 7(b), Opt-S has less downcoding cost than Opt-R for scaling operation p_5 , and the reason is that Opt-S can exploit the local parity block of the compact LRC to calculate one remaining local parity block, therefore saving one network transfer in each upcoding unit (Section III-D). However, for scaling operations p_2, p_4 and p_6 , Opt-R can regenerate each local parity block within each cluster locally, and thus has less downcoding cost (zero) than Opt-S.
- From Figure 7(d), the repair cost of the compact LRC of Opt-S equals that of Opt-R (because the layout of the compact LRC of Opt-S is the same as that of Opt-R (Figure 4)), and both are much smaller than that of Flat.
- When g increases, the scaling and repair costs of both Opt-S and Opt-R decrease, while these costs of Flat keep constant. Hence, the improvements of Opt-S and Opt-R over Flat become greater with larger g . For example, Opt-S reduces the repair cost of the fast LRC of Flat by 67% for scaling operation p_2 (i.e., $g = 2$), while the reduction becomes 83% for scaling operation p_3 (i.e., $g = 4$).

B. Testbed Experiments

We implement Opt-S, Opt-R, and Flat in a distributed storage system prototype, and conduct testbed experiments to understand their scaling and repair performance. In our prototype, repair and scaling are both implemented as two phase processes, where we first aggregate relevant data within a cluster, and then send the aggregated data across cluster to the destination cluster. Our prototype is written in C++ and

implemented with a Coordinator (CN) and multiple Datanodes (DN). The CN sends commands, while the DNs receive commands and execute the actual data read, write, and transfer independently and in parallel, and finally reply acks to the CN.

Setup. Our testbed comprises 22 physical nodes, each of which runs Ubuntu 16.04.5 LTS with a quad-core 3.40 GHz Intel Core i5-3570, 16 GB RAM, and a Seagate ST1000DM003 7200 RPM 1 TB SATA hard disk. Each node achieves 170 MBps of disk read bandwidth and 130 MBps of write bandwidth, and 1 Gbps of network bandwidth (measured by iperf). We deploy the CN in one node, and the DNs on 20 nodes. We configure the DNs into different clusters according to the scaling operations and the placement strategy. We also configure a dedicated node to act as a network core, such that any cross-cluster traffic must traverse the network core. We use the Wonder Shaper tool [1] to control the outgoing bandwidth of the network core.

Methodology. We assume the following default configurations. We adopt the scaling operation p_2 (i.e., (12, 6, 2, 7) to (12, 2, 2, 7)). We configure the block size as 64 MB, the packet size as 1 MB (packet is the unit for network transfer), and the cross-cluster bandwidth as 100 Mbps (such that the ratio of inner-cluster bandwidth to cross-cluster bandwidth is 10:1). We may vary some of the settings in our experiments. We measure the repair time per block and the scaling time per stripe. The results of each experiment are averaged over five runs.

Experiment 1 (Performance for different scaling operations). We first evaluate the performance for different scaling operations. We consider three sets of scaling operations, i.e., p_0, p_1 and p_2 . We fix the block size as 64 MB and the cross-cluster bandwidth as 100 Mbps, and then compare the scaling time and the repair time. Figure 8 shows the results.

According to our analysis, the upcoding cost of Opt-S is zero, while that of Flat and Opt-R increases as $l - l'$ increases. From Figure 8(a), the experimental results comply with the theoretics. As $l - l'$ grows, the upcoding time of Opt-S keeps fairly stable while that of Flat and Opt-R increases. Thus, Opt-S will achieve more gain for larger $l - l'$. Overall, Opt-S reduces the upcoding time of Flat and Opt-R by 80.6% and 80.2%, 87.6% and 87.2%, and 91.0% and 91.0%, for p_0, p_1 and p_2 , respectively. The upcoding time of Opt-R is similar to that of Flat, which confirms to the numerical results.

The repair cost of Opt-R is zero while that of Flat keeps constant as b keeps unchanged. The repair cost of Opt-S (i.e., $1 - \frac{\theta}{\delta}$) grows as $\delta = \frac{l}{l'}$ increases. From Figure 8(c), the repair time of Opt-R and Flat keeps stable while that of Opt-S grows slightly. However, Opt-R always has the smallest repair time. Overall, Opt-R reduces the repair time of Flat and Opt-S by 85.5% and 56.8%, 84.2% and 55.8%, and 84.0% and 61.5%, for scaling operations p_0, p_1 and p_2 , respectively.

From Figure 8(b), both Opt-S and Opt-R show better downcoding time performance compared to Flat. For example, for scaling operation p_2 , Opt-S and Opt-R reduce the downcoding time of Flat by 80.7% and 97.3%, respectively.

Finally, from Figure 8(d), the repair time of the compact LRC of Opt-S is almost identical to that of Opt-R, and both show

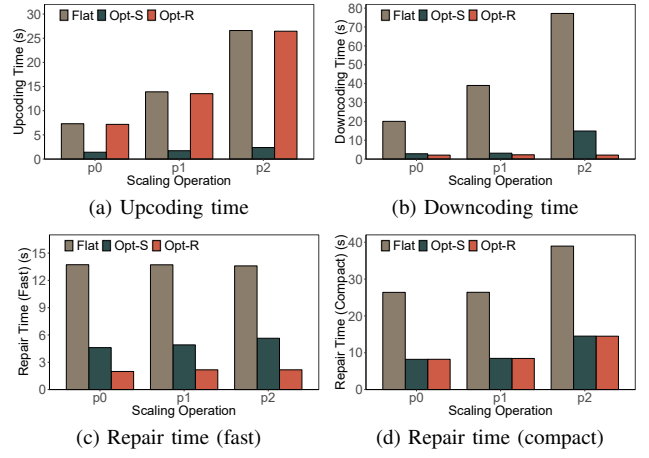


Fig. 8. Experiment 1: Comparison of the scaling and repair time.

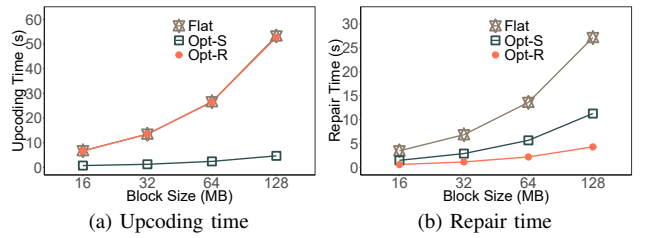


Fig. 9. Experiment 2: Impact of the block size.

improvements over Flat. For example, for scaling operation p_2 , Opt-S and Opt-R reduce the repair time of Flat by 62.8% and 62.8%, respectively.

Experiment 2 (Impact of block size). We now evaluate the impact of the block size, varied from 16 MB to 128 MB. We test the default scaling operation (i.e., p_2) and fix the cross-cluster bandwidth as 100 Mbps. We only show the results for the upcoding time and the repair time of the fast LRC in the following experiments. Figure 9 shows the results. We can see that the upcoding time and the repair time increases with a larger block size, and Opt-S and Opt-R constantly outperform Flat. For example, Opt-S reduces the upcoding time of Flat from 89.6%-91.3%, and Opt-R reduces the repair time of Flat from 83.2%-84.3%, across all block sizes.

Experiment 3 (Impact of bandwidth). We now study the impact of cross-cluster bandwidth. Here, we adopt p_2 and fix the block size as 64 MB, and then vary the cross-cluster bandwidth from 50 Mbps to 400 Mbps (the ratio of inner-cluster bandwidth to cross-cluster bandwidth is 20:1-2.5:1). Figure 10 shows the results. As expected, the upcoding and repair time decreases with larger bandwidth. Besides, in Figure 10(a), Opt-S reduces the upcoding time of Flat by 95.2%, 91.0%, 84.1% and 74.7% when the bandwidth is 50 Mbps, 100 Mbps, 200 Mbps and 400 Mbps, respectively. In Figure 10(b), Opt-R reduces the repair time of Flat by 91.5%, 84.0%, 72.6% and 56.6% when the bandwidth changes. These indicate that the improvements of Opt-S and Opt-R over Flat are greater with more scarce cross-cluster bandwidth.

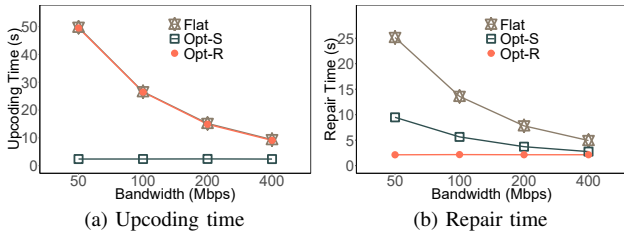


Fig. 10. Experiment 3: Impact of the cross-cluster bandwidth.

V. RELATED WORK

There has been extensive work on the repair performance of LRC in the literature. Theoretical studies on LRC (e.g., [8], [23]–[25]) focus on the relationship between the optimal minimum Hamming distance and the theoretical repair cost, and design explicit LRC constructions. LRC is also implemented and evaluated in Azure [11] and Facebook [21]. Kolosov *et al.* [14] study the trade-offs of different LRC constructions between storage overheads and repair costs. In contrast, our work mainly focuses on the trade-off between the repair and scaling costs in hierarchical clustered storage systems.

Several studies address the scaling problem on the change of erasure coding configurations. Some studies (e.g., [4], [12], [30], [31], [33], [34]) explore the efficient data redistribution when a storage cluster expands and adds new storage nodes, while others (e.g., [6], [15], [28], [29]) address the transition from replication to erasure coding. HeART [13] suggests the change of redundancy settings to balance between storage overhead and reliability. To efficiently adapt to workload changes, HACFS [32] proposes efficient transitions between two erasure codes to balance between storage overhead and access performance. Our work is motivated by the scenario in HACFS, and presents a formal repair-scaling trade-off analysis.

A number of studies consider the deployment of erasure coding in clustered storage that spans multiple geographic regions [3], [20]. In particular, some studies focus on minimizing the cross-cluster bandwidth during repair [9], [10], [16], [18], [22]. Our work is the first one that addresses scaling performance in clustered storage.

VI. CONCLUSION

We investigate the optimal trade-off between the repair and scaling performance of LRC in clustered storage systems. We design placement strategies that operate along the optimal repair-scaling trade-off curve subject to the single-cluster fault tolerance constraint. Both numerical studies and testbed experiments validate the efficiency of our placement strategies. The source code of our implementation is available at <http://adslab.cse.cuhk.edu.hk/software/lrctradeoff>.

Acknowledgements: This work is supported by the Research Grants Council of Hong Kong (GRF 14216316 and CRF C7036-15G), the National Natural Science Foundation of China (61602120), and the Fujian Provincial Natural Science Foundation (2017J05102).

APPENDIX

Proof of Theorem 1. Suppose that the upcoding cost of a placement is u , and the upcoding cost induced by the i -th ($0 \leq i \leq l' - 1$) upcoding unit is u_i ($\sum_{i=0}^{l'-1} u_i = u$). For an upcoding unit, there must be a *core cluster* that will collect and encode the δ local parity blocks. There must also be some different clusters that hold local parity blocks. We call each of such clusters a *remote cluster*. Our proof is organized as follows.

(i) We first prove that the number of remote clusters for the i -th upcoding unit is exactly u_i . During upcoding, the core cluster will retrieve an XOR sum of the local parity blocks in each remote cluster. If the number of remote clusters is not u_i , then the upcoding cost will be different from u_i .

(ii) We next prove that a core cluster also holds local parity blocks. Otherwise, we can choose one of the remote clusters as the new core cluster. During upcoding, the local parity blocks in the new core cluster can be retrieved locally, so the upcoding cost for the i -th upcoding unit will be $u_i - 1$ rather than u_i .

(iii) We then prove that we can collocate (at most) θ local data sets into a cluster that holds local parity blocks if the corresponding local parity blocks of these local data sets are also in this cluster. Suppose that a cluster holds r local parity blocks, and we can collocate (at most) w local data sets whose local parity blocks are included in the r ones. The number of blocks (i.e., $w \times b + r$), which span r local groups, cannot exceed $g + r$ (Lemma 1). Thus, $w = \lfloor \frac{g}{b} \rfloor = \theta$, and the cost for repairing any data block in these θ local data sets is zero.

(iv) We now show that the core cluster and remote clusters of one upcoding unit should be different from those of another upcoding unit so as to minimize the repair cost. If cluster A of one upcoding unit is the same as cluster B of another upcoding unit, then according to (iii), we can put (at most) θ local data sets into A (or B). If A is different from B , then we can put θ local data sets into A and another θ local data sets into B . Hence, if the core cluster and remote clusters for each upcoding unit are different, then we can collocate more local data sets to be together with their local parity blocks, and so the repair cost of more data blocks is zero.

(v) We now prove that the maximum number of data blocks whose repair cost is zero is $(l' + u) \times \theta \times b$. According to (iii), we can put (at most) θ local data sets into a cluster that holds local parity blocks, so as to make the cost for repairing any data block in these θ local data sets as zero. According to (i), (ii) and (iv), we have (at most) l' core clusters and u remote clusters that hold local parity blocks. Thus, the maximum number of data blocks whose repair cost is zero is $(l' + u) \times \theta \times b$.

(vi) Finally, we prove that we cannot have a placement whose repair cost is less than $1 - \frac{\theta}{\delta} - \frac{u \times \theta}{l}$ subject to the upcoding cost of u . Otherwise, since the repair cost is less than one, there must be some data blocks whose repair cost is zero. We assume that the number of data blocks whose repair cost is zero is α , so the number of data blocks whose repair cost is at least one is $k - \alpha$. Therefore, $\frac{k - \alpha}{k} < 1 - \frac{\theta}{\delta} - \frac{u \times \theta}{l}$, implying $\alpha > (l' + u) \times \theta \times b$. However, this is conflicting with our proof in (v). Thus, we complete our proof.

REFERENCES

- [1] The Wonder Shaper 1.4. <https://github.com/magnifico/wondershaper>.
- [2] F. Ahmad, S. T. Chakradhar, A. Raghunathan, and T. Vijaykumar. ShuffleWatcher: Shuffle-aware Scheduling in Multi-tenant Mapreduce Clusters. In *Proc. of USENIX ATC*, 2014.
- [3] Y. L. Chen, S. Mu, J. Li, C. Huang, J. Li, A. Ogus, and D. Phillips. Giza: Erasure Coding Objects across Global Data Centers. In *Proc. of USENIX ATC*, 2017.
- [4] L. Cheng, Y. Hu, and P. P. Lee. Coupling Decentralized Key-Value Stores with Erasure Coding. In *Proc. of ACM SoCC*, 2019.
- [5] M. Chowdhury, S. Kandula, and I. Stoica. Leveraging Endpoint Flexibility in Data-Intensive Clusters. In *Proc. of ACM SIGCOMM*, 2013.
- [6] B. Fan, W. Tantisiriroj, L. Xiao, and G. Gibson. DiskReduce: RAID for Data-Intensive Scalable Computing. In *Proc. of ACM PDSW*, 2009.
- [7] D. Ford, F. Labelle, F. Popovici, M. Stokely, V.-A. Truong, L. Barroso, C. Grimes, and S. Quinlan. Availability in Globally Distributed Storage Systems. In *Proc. of USENIX OSDI*, 2010.
- [8] P. Gopalan, C. Huang, H. Simitci, and S. Yekhanin. On the Locality of Codeword Symbols. *IEEE Trans. on Information Theory*, 58(11):6925–6934, 2012.
- [9] H. Hou, P. P. Lee, K. W. Shum, and Y. Hu. Rack-Aware Regenerating Codes for Data Centers. *IEEE Trans. on Information Theory*, 65(8):4730–4745, Aug 2019.
- [10] Y. Hu, X. Li, M. Zhang, P. P. Lee, X. Zhang, P. Zhou, and D. Feng. Optimal Repair Layering for Erasure-Coded Data Centers: From Theory to Practice. *ACM Trans. on Storage*, 13(4):33, 2017.
- [11] C. Huang, H. Simitci, Y. Xu, A. Ogus, B. Calder, P. Gopalan, J. Li, and S. Yekhanin. Erasure Coding in Windows Azure Storage. In *Proc. of USENIX ATC*, 2012.
- [12] J. Huang, X. Liang, X. Qin, P. Xie, and C. Xie. Scale-RS: An Efficient Scaling Scheme for RS-Coded Storage Clusters. *IEEE Trans. on Parallel and Distributed Systems*, 26(6):1704–1717, 2014.
- [13] S. Kadekodi, K. Rashmi, and G. R. Ganger. Cluster storage systems gotta have HeART: improving storage efficiency by exploiting disk-reliability heterogeneity. In *Proc. of USENIX FAST*, 2019.
- [14] O. Kolosov, G. Yadgar, M. Liram, I. Tamo, and A. Barg. On Fault Tolerance, Locality, and Optimality in Locally Repairable Codes. In *Proc. of USENIX ATC*, 2018.
- [15] R. Li, Y. Hu, and P. P. Lee. Enabling Efficient and Reliable Transition from Replication to Erasure Coding for Clustered File Systems. *IEEE Trans. on Parallel and Distributed Systems*, 28(9):2500–2513, 2017.
- [16] X. Li, R. Li, P. P. Lee, and Y. Hu. OpenEC: Toward Unified and Configurable Erasure Coding Management in Distributed Storage Systems. In *Proc. of USENIX FAST*, 2019.
- [17] S. Muralidhar, W. Lloyd, S. Roy, C. Hill, E. Lin, W. Liu, S. Pan, S. Shankar, V. Sivakumar, L. Tang, et al. f4: Facebook’s Warm BLOB Storage System. In *Proc. of USENIX OSDI*, 2014.
- [18] N. Prakash, V. Abdrashitov, and M. Médard. The Storage versus Repair-Bandwidth Trade-off for Clustered Storage Systems. *IEEE Trans. on Information Theory*, 64(8):5783–5805, Aug 2018.
- [19] K. Rashmi, N. B. Shah, D. Gu, H. Kuang, D. Borthakur, and K. Ramchandran. A Solution to the Network Challenges of Data Recovery in Erasure-coded Distributed Storage Systems: A Study on the Facebook Warehouse Cluster. In *Proc. of USENIX HotStorage*, 2013.
- [20] J. K. Resch and J. S. Plank. AONT-RS: Blending Security and Performance in Dispersed Storage Systems. In *Proc. of USENIX FAST*, 2011.
- [21] M. Sathiamoorthy, M. Asteris, D. Papailiopoulos, A. G. Dimakis, R. Vadali, S. Chen, and D. Borthakur. XORing Elephants: Novel Erasure Codes for Big Data. In *Proc. of VLDB Endowment*, pages 325–336, 2013.
- [22] Z. Shen, J. Shu, and P. P. Lee. Reconsidering Single Failure Recovery in Clustered File Systems. In *Proc. of IEEE/IFIP DSN*, 2016.
- [23] N. Silberstein, A. S. Rawat, O. O. Koyluoglu, and S. Vishwanath. Optimal Locally Repairable Codes via Rank-Metric Codes. In *Proc. of IEEE International Symposium on Information Theory*, 2013.
- [24] I. Tamo and A. Barg. A Family of Optimal Locally Recoverable Codes. *IEEE Trans. on Information Theory*, 60(8):4661–4676, 2014.
- [25] I. Tamo, D. S. Papailiopoulos, and A. G. Dimakis. Optimal Locally Repairable Codes and Connections to Matroid Theory. *IEEE Trans. on Information Theory*, 62(12):6661–6671, 2016.
- [26] A. Vulimiri, C. Curino, P. B. Godfrey, T. Jungblut, J. Padhye, and G. Varghese. Global Analytics in the Face of Bandwidth and Regulatory Constraints. In *Proc. of USENIX NSDI*, 2015.
- [27] H. Weatherspoon and J. D. Kubiatowicz. Erasure Coding Vs. Replication: A Quantitative Comparison. In *Proc. of IPTPS*, 2002.
- [28] S. Wei, Y. Li, Y. Xu, and S. Wu. DSC: Dynamic Stripe Construction for Synchronous Encoding in Clustered File System. In *Proc. of IEEE INFOCOM*, 2017.
- [29] J. Wilkes, R. Golding, C. Staelin, and T. Sullivan. The HP AutoRAID Hierarchical Storage System. *ACM Trans. on Computer Systems*, 14(1):108–136, 1996.
- [30] C. Wu and X. He. GSR: A Global Stripe-Based Redistribution Approach to Accelerate RAID-5 Scaling. In *Proc. of IEEE ICPP*, 2012.
- [31] S. Wu, Y. Xu, Y. Li, and Z. Yang. I/O-Efficient Scaling Schemes for Distributed Storage Systems with CRS Codes. *IEEE Trans. on Parallel and Distributed Systems*, 27(9):2639–2652, 2016.
- [32] M. Xia, M. Saxena, M. Blaum, and D. A. Pease. A Tale of Two Erasure Codes in HDFS. In *Proc. of USENIX FAST*, 2015.
- [33] X. Zhang, Y. Hu, P. P. Lee, and P. Zhou. Toward Optimal Storage Scaling via Network Coding: From Theory to Practice. In *Proc. of IEEE INFOCOM*, 2018.
- [34] W. Zheng and G. Zhang. FastScale: Accelerate RAID Scaling by Minimizing Data Migration. In *Proc. of USENIX FAST*, 2011.