

Achieving Load-Balanced, Redundancy-Free Cluster Caching with Selective Partition

Yinghao Yu, *Student Member, IEEE*, Wei Wang, *Member, IEEE*, Renfei Huang, *Student Member, IEEE*, Jun Zhang, *Senior Member, IEEE*, and Khaled B. Letaief, *Fellow, IEEE*

Abstract—Data-intensive clusters increasingly rely on in-memory storages to improve I/O performance. However, the routinely observed file popularity skew and load imbalance create hot spots, which significantly degrade the benefits of in-memory caching. Common approaches to tame load imbalance include copying multiple replicas of hot files and creating parity chunks using storage codes. Yet, these techniques either suffer from high memory overhead due to cache redundancy or incur non-trivial encoding/decoding complexity. In this paper, we propose an effective approach to achieve load balancing without cache redundancy or encoding/decoding overhead. Our solution, termed SP-Cache, *selectively partitions* files based on the loads they contribute and evenly caches those partitions across the cluster. We develop an efficient algorithm to determine the optimal number of partitions for a hot file—too few partitions are incapable of mitigating hot spots, while too many are susceptible to stragglers. We have implemented SP-Cache atop Alluxio, a popular in-memory distributed storage system, and evaluated its performance through EC2 deployment and trace-driven simulations. SP-Cache can quickly react to the changing load by dynamically re-balancing cache servers. Compared to the state-of-the-art solution, SP-Cache reduces the file access latency by up to 40% in both the mean and the tail, using 40% less memory.

Index Terms—Cloud computing, cluster caching systems, load balancing, selective partition

1 INTRODUCTION

Today’s data-parallel clusters employ *in-memory* storage solutions for high-performance data analytics [2]–[7]. By caching data objects in memory, I/O-intensive applications can gain order-of-magnitude performance improvement over traditional on-disk solutions [2], [4], [5].

However, one plaguing problem faced by in-memory solutions is the severe *load imbalance* across cache servers. In production clusters, data objects typically have the *heavily skewed popularity*—meaning, a small number of hot files account for a large fraction of data accesses [8]–[10]. The cache servers storing hot files hence turn into *hot spots*. This problem is further aggravated by the *network load imbalance*. It was reported in a Facebook cluster that the most heavily loaded links have over $4.5\times$ higher utilization than the average for more than 50% of the time [8]. The routinely observed hot spots, along with the network load imbalance, result in a significant degradation of I/O performance that could even cancel out the performance advantage of in-memory caching (Sec. 2).

Therefore, maintaining load balance across cache servers is the key to improving the performance of cluster caches. State-of-the-art solutions in this regard include *selective replication* [9] and *erasure coding* [8], both of which resort to *redundant caching* to mitigate hot spots.

Selective replication creates multiple replicas for files based on their popularity: the more popular a file is, the more replicas it has. File access requests can then be distributed to multiple servers containing those replicas, hence mitigating the load on the hot spots. However, replication results in high memory overhead as hot files are usually of large sizes [8]–[10]. Given the limited memory space, selective replication does not perform well in cluster caches [8], [11] (Sec. 3.1).

Erasure coding comes as an alternative solution to achieve load balancing with reduced memory overhead [8]. In particular, a (k, n) erasure code divides a file into k partitions and generates $n - k$ *parity partitions*. Any k of the n partitions are sufficient to decode the original file. This results in better load balancing as the load of read requests are spread across multiple servers, and the memory overhead is usually much smaller than that of replication. However, erasure coding incurs salient encoding/decoding overhead. In fact, even with the highly optimized implementation [12], the computational overhead can still delay the I/O requests by 30% on average [8].

In this paper, we propose a different approach that achieves load balancing in cluster caches *without* memory redundancy or encoding/decoding overhead. Our approach, which we call *selective partition*, divides files into multiple partitions based on their popularity: the more popular a file is, the more partitions it is split into. File partitions are *randomly cached* by servers across the cluster. The benefits of this approach are three-fold. First, it *evenly spreads* the load of read requests across cache servers, leading to improved

- Y. Yu and K. B. Letaief are with the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong (e-mail: {yyuau, eekhaled}@ust.hk). K. B. Letaief is also an Adjunct Distinguished Scientist at Peng Cheng Laboratory in Shenzhen.
- W. Wang and R. Huang are with the Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong (e-mail: {weiw, rhuangan}@ust.hk).
- J. Zhang is with the Department of Electronic and Information Engineering, Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong (e-mail: jun-eie.zhang@polyu.edu.hk).
- Part of this paper has appeared in [1]. This new version contains substantial revision with theoretic analysis, improved implementations and more comprehensive evaluations.

load balancing. Second, it increases the *read parallelism* of hot files, which, in turn, improves the I/O performance. Third, simply splitting files into partitions adds *no storage redundancy*, nor does it incur the computational overhead for encoding/decoding.

However, it remains a challenge to *judiciously* determine how many partitions a file should be split into. On one hand, too few partitions are insufficient to spread the load of read requests, making it incapable of mitigating hot spots. On the other hand, reading too many partitions from across servers adds the risk of being slowed down by *stragglers*.

To address this challenge, we model selective partition as a *fork-join* queuing system [13], [14] and establish an *upper-bound analysis* to quantify the mean latency in reads. We show that the optimal number of partitions can be efficiently obtained by solving a convex optimization problem. Based on this result, we design SP-Cache, a *load-balanced, redundancy-free* cluster caching scheme that optimally splits files to minimize the mean latency while mitigating the impact of stragglers. We show that SP-Cache improves load balancing by a factor of $O(L_{\max})$ compared to the state-of-the-art solution called EC-Cache [8], where L_{\max} measures the load of the *hottest* file.

We have implemented SP-Cache atop Alluxio [2], [15], a popular in-memory distributed storage that can be used as the caching layer on top of disk-based cloud object stores (e.g., Amazon S3 [16] and Azure Storage [17]) and compute-located cluster file systems (e.g., HDFS [18] and Gluster [19]). We evaluated SP-Cache through both EC2 [20] deployment and trace-driven simulations. Experimental results show that despite the presence of intensive stragglers, SP-Cache reduces the mean and the tail (95th percentile) read latency by up to 40% compared to EC-Cache [8]. Owing to its redundancy-free nature, SP-Cache achieves all these benefits with 40% less memory footprint than EC-Cache.

Part of this work has appeared in [1]. Compared to the conference version [1], we have made several substantial improvements in this paper:

- 1) An efficient implementation of SP-Cache with parallel repartition scheme (Sec. 6.2), which enables SP-Cache to quickly react to the changing file popularity with a speedup of two orders of magnitude over the previous implementation [1] (Sec. 7.4);
- 2) A detailed description of the configuration of SP-Cache (Algorithm 1);
- 3) More comprehensive evaluations of the repartition overhead (Sec. 7.4), the write performance (Sec. 7.8), and the coefficient of variation in read latency (Secs. 2.2, 3.1, and 4.1);
- 4) Expanded discussions of future extensions (Sec. 8) and related works (Sec. 9).

2 BACKGROUND AND MOTIVATION

In this section, we briefly survey the cluster caching systems and motivate the need to achieve load balancing therein.

2.1 Cluster Caching

Due to the recent technological advances in datacenter fabrics [21] and the emergence of new high-speed network

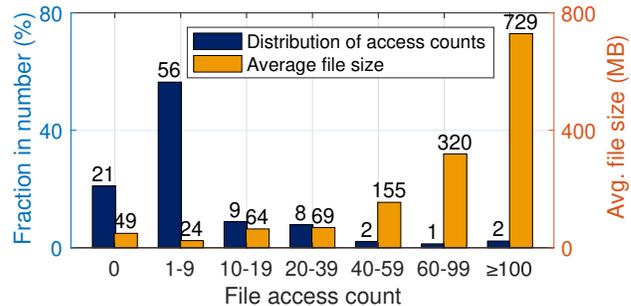


Fig. 1: Distribution of access counts (blue) and average file size (orange) observed in a Yahoo! cluster [33].

appliances [22]–[24], the gap between network bandwidth and storage I/O bandwidth is rapidly narrowing [25]–[28]. Consequently, the performance bottleneck of cloud systems is shifting from network to storage I/O. Prior work has shown that accessing data from the local hard disk provides no salient benefits over remote reads [29], [30]. As disk locality becomes irrelevant, cloud object stores, such as Amazon S3 [16], Windows Azure Storage [17], and OpenStack Swift [31], gradually replace compute-located storages— notably HDFS [18]—as the primary storage solutions for data-intensive applications.

However, cloud object stores remain bottlenecked on disk I/O [8], as reading from disk is at least two orders of magnitude slower than reading from memory. In light of this problem, cluster caching systems, such as Alluxio [15], Memcached [6], and Redis [32], are increasingly deployed in front of cloud object stores to provide low-latency data access at memory speed. In this paper, we primarily target the storage-side caching to improve the I/O performance. Our solution can also be applied to compute-located file systems, such as HDFS [18], provided that high-speed networks are available.

2.2 Load Imbalance and Its Impact

A plaguing problem faced by cluster caching is the routinely observed *load imbalance* across cache servers. We show through experiments that severe load imbalance results in significant I/O latencies, marginalizing the performance benefits provided by cluster caching.

Load Imbalance Prior works [8], [9] have identified two sources of load imbalance in production clusters: the *skewed file popularity* and the *imbalanced network traffic*.

It has been widely observed in datacenters that file (data object) popularity is heavily skewed and usually follows a Zipf-like distribution [2], [8]–[10]. That is, a large fraction of data access requests are contributed by only a small number of hot files. Fig. 1 depicts the distribution of file access counts and size in the Yahoo! cluster trace [33]. The trace contains the collective statistics of data accesses to over 40 million files in a period of two months. We observe that the majority of files (~ 78%) have cold data that has rarely been accessed (< 10 times). Only 2% are hot with high access counts (≥ 100). These files are usually much larger (15–30×) than the cold ones. Consequently, cache servers containing these files are easily overloaded given their large sizes and high popularity.

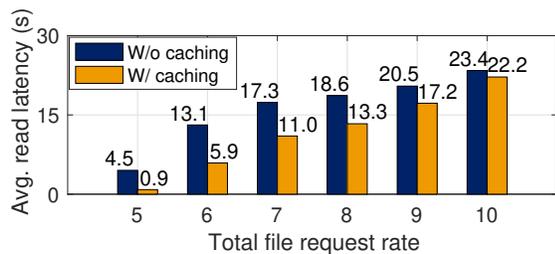


Fig. 2: The average read latencies with and without caching as the load of read requests increases.

This problem is further aggravated in the presence of network load imbalance, which is prevalent in production datacenters [8], [34]–[36]. For example, a recent study [8] measured the ratio of the maximum and the average utilizations across all up- and down-links in a Facebook cluster. The result shows that the ratio stays above $4.5\times$ more than half of the time, suggesting a severe imbalance.

Impact of Load Imbalance The skew in file popularity, together with the imbalanced network traffic, create hot spots among cache servers. To illustrate how these overloaded machines may impair the system’s I/O performance, we stress-tested a small cluster of cache servers.

Setup: We deployed Alluxio [15]—a popular in-memory distributed storage—on a 30-node Amazon EC2 [20] cluster. The nodes we used are *m4.large* instances, each with a dual-core processor, 8 GB memory, and 0.8 Gbps network bandwidth. The cluster is used to cache 50 files (40 MB each). We launched another 20 *m4.large* instances as *clients*. Each client submits the file read requests to the Alluxio cluster as a Poisson process with a rate from 0.25 to 0.5 requests per second. Therefore, the aggregated access rates are 5–10 requests per second. We created imbalanced load with skewed file popularity following a Zipf distribution with exponent 1.1 (i.e., high skewness).

Diminishing benefits of caching: We ran two experiments. In the first experiment, all files were cached in memory; in the second experiment, we disabled cluster caching and spilled files to the local hard disk. For each experiment, we measured the mean read latency under various request rates and depict the results in Fig. 2. The coefficients of variance (CV) are shown in Table 1. Note that having CV greater than 1 indicates high variance in read latency, i.e., severe hot spot effects.

When the cluster is less loaded (5 requests per second), in-memory caching provides salient benefits, improving the mean read latency by $5\times$. However, as the load ramps up, the hot spots among cache servers become more pronounced, and the benefits of caching quickly diminish. We make similar observations for the CV in read latencies. Under the skewed popularities, the CV is consistently higher than 1, suggesting high variance due to the presence of hot spots. In fact, with request rate greater than 9, the read latency is dominated by the network congestion on hot-spot servers, and in-memory caching becomes *irrelevant*.

Therefore, there is a pressing need to achieve load balancing across servers. We next review existing techniques and show that they all fall short in minimizing latency.

TABLE 1: The coefficient of variation (CV) of the read latencies with and without caching as the load of read requests increases.

Request rate	5	6	7	8	9	10
W/o caching	1.67	1.70	1.64	1.74	1.79	1.78
W/ caching	1.29	1.41	1.59	2.08	1.83	1.83

TABLE 2: The coefficient of variation (CV) of the read latencies with different replica numbers of the top 10% popular files.

Replication #	1	2	3	4	5
CV	1.29	1.25	1.22	0.61	0.64

3 INEFFICIENCY OF EXISTING SOLUTIONS

Prior art resorts to *redundant caching* to achieve load balancing, either by *copying multiple replicas* of hot files—known as *selective replication* [9], [37]—or by creating *coded partitions* of data objects, e.g., EC-Cache [8]. However, these techniques enforce an *unpleasant trade-off* between load balancing and cache efficiency. On one hand, caching more replicas (coded partitions) helps mitigate hot spots as the load of read requests can be spread to more servers. On the other hand, the overhead in memory and/or computation due to redundant caching harms efficiency.

3.1 Selective Replication

Selective replication replicates files based on their popularity [9], [37], i.e., the more popular a file is, the more replicas are copied across servers. A file read request is then randomly served by one of the servers containing the replica of that file. In this way, the load of read requests are evenly distributed, leading to improved load balancing.

While selective replication is proven effective for disk-based storage [9], it does not perform well for cluster caching [8], [11], as replication incurs high memory overhead. To illustrate this problem, we deployed an Alluxio cluster following the settings described in Sec. 2.2, where the top 10% popular files were copied to multiple replicas. The aggregated request rate is set to 6. We gradually increased the number of replicas and examined how the mean latency in reads can be improved at the expense of increased memory overhead. Fig. 3 and Table 2 depict the average latencies and coefficients of variances, respectively.

We observe a *linear growth* of memory overhead in exchange for only a *sublinear improvement* in read latency. Moreover, we need a replication factor of 4 to effectively suppress the coefficient of variation to be lower than 1. In other words, to mitigate the hot spots, we need 3 extra copies for each of the top 10% popular files. Given that in-memory caches remain a constrained resource in production clusters and the fact that popular files are usually of large sizes (Fig. 1), selective replication often results in poor cache efficiency with very low hit ratio (more in Sec. 7.7).

3.2 Erasure Coding

State-of-the-art solutions employ *erasure coding* [38], [39] to load-balance cache servers without incurring high memory overhead. In particular, a (k, n) erasure coding scheme

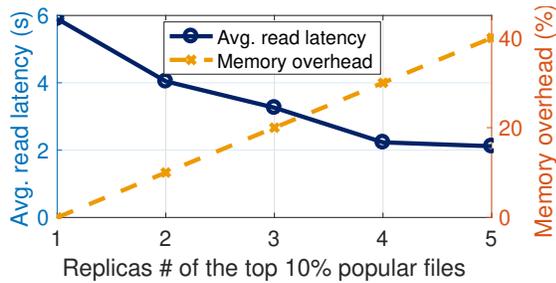


Fig. 3: Average read latency and cache cost in percentage with different replica numbers of the top 10% popular files.

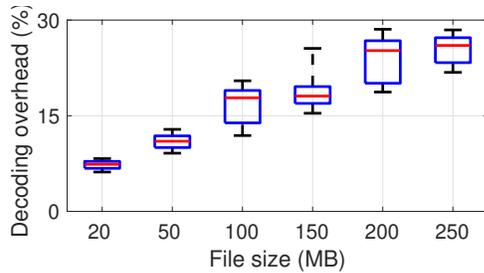


Fig. 4: Decoding overhead in EC-Cache. Boxes depict the 25th, 50th, and 75th percentiles. Whiskers depict the 5th and 95th percentiles.

evenly splits a file into k partitions. It then computes $n - k$ parity partitions of the same size. The original file can be decoded using any k out of the n partitions, allowing the load of its read requests to be spread to n servers. The memory overhead is $(n - k)/k$, which is lower than that of selective replication (at least $1\times$) in practical settings. An efficient implementation of this approach goes to EC-Cache [8] which *late-binds* partitions during reads to mitigate stragglers. That is, instead of reading exactly k partitions, EC-Cache randomly fetches $k + 1$ partitions and waits for any k partitions to complete reading. EC-Cache significantly outperforms selective replication in both the median and tail read latencies [8].

However, EC-Cache requires non-trivial decoding (encoding) overhead during reads (writes). Even with a highly optimized coding scheme [40] and implementation [12], the decoding overhead may still delay the read requests by up to 30% [8]. To verify this result, we ran EC-Cache in an Amazon EC2 cluster with 30 r3.2xLarge memory-optimized instances, each having 61 GB memory and 8 cores. Following [8], we used a (10, 14) coding scheme (i.e., memory overhead 40%) to cache files of various sizes. We launched an EC-Cache client submitting file read requests and measured the incurred decoding overhead, i.e., the decoding time *normalized* by the read latency. Fig. 4 depicts the results as a box plot. We observed more prominent decoding overhead with large files. Notably, for files greater than 100 MB which account for most of the file accesses in production clusters (Fig. 1), the decoding overhead consistently stays above 15%. We stress that this result is measured in the presence of a less advanced network, where we observed 1 Gbps bandwidth between instances. We expect the read latency dominated by the decoding overhead with high-

TABLE 3: The coefficient of variation (CV) of read latencies using simple partition in a 30-node cluster, with and without stragglers.

Partition #	3	9	15	21	27
W/o stragglers	1.02	0.75	0.55	0.44	0.48
W/ stragglers	1.03	1.10	1.05	1.17	1.35

speed networks (≥ 40 Gbps bisection bandwidth).

To sum up, existing load balancing solutions either suffer from high cache redundancy or incur non-trivial decoding/encoding overhead—either way, the I/O performance is impaired.

4 LOAD BALANCING WITH SIMPLE PARTITION

In this section, we consider a simple, yet effective load-balancing technique which uniformly splits files into multiple partitions so as to spread the I/O load. We explore the potential benefits as well as the problems it causes.

4.1 Simple Partition and Potential Benefits

We learn from EC-Cache [8] that dividing files into smaller partitions improves load balancing, yet the presence of parity partitions necessitates the decoding overhead. A simple fix is to split files *without creating coded partitions*. Intuitively, simple partition provides two benefits over EC-Cache. *First*, it requires no overhead for decoding/encoding. *Second*, it adds no storage redundancy, attaining the highest possible cache utilization. Simple partition also retains two benefits provided by EC-Cache. *First*, it mitigates hot spots under skewed popularity by spreading the load of read requests to multiple partitions. *Second*, it provides opportunities for read/write parallelism, which, in turn, speeds up the I/O of large files.

To validate these potential benefits, we configured EC-Cache in a “coding-free” mode using a (k, k) coding scheme (i.e., no parity partition). Specifically, we evenly split a file into k partitions and randomly cached them across the cluster. No two partitions of a file were placed on the same server. We re-ran the experiments in Sec. 2.2 using simple partition. For the purpose of stress-testing, we configured the aggregated file request rate to 10 from all clients. Note that at such a high rate, the average read latency would have stretched over 20 s *without load balancing* (Fig. 2). We study how simple partition can speed up I/O with increased read parallelism k . The results are depicted as a solid line in Fig. 5. The average read latencies drop to 1-1.3 s, suggesting 17-22 \times improvement over stock Alluxio without partition. Table 3 presents the coefficient of variance, which is reduced to 0.75 with 9 partitions and below 0.5 with more partitions. We stress that these improvements are achieved without any decoding overhead or cache redundancy. In contrast, the replication scheme studied in Sec. 3.1 is only able to attain the average latency of 2 s and a CV of 0.61 with 1.4 \times memory footprint (Fig. 3) in the presence of even lighter load (i.e., 6 requests per second).

4.2 Problems of Simple Partition

However, simple partition is not without problems. First, it *uniformly* divides each file into k partitions, irrespective of its

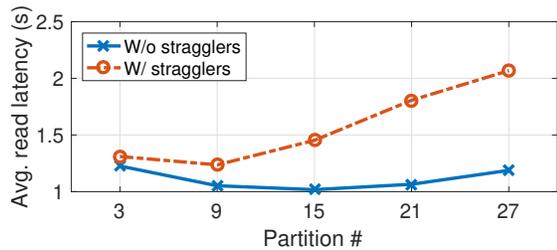


Fig. 5: Average read latency using simple partition in a 30-node cluster, with and without stragglers.

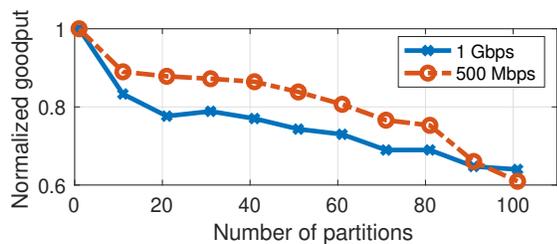


Fig. 6: Normalized goodput with different number of partitions using 1 Gbps and 500 Mbps network, respectively.

size and popularity. This is *unnecessary* and *inefficient*. For less popular files, which usually dominate in population [9], [10], spreading their load provides marginal benefits in improving load balancing. Rather, the increased read parallelism may result in salient networking overhead due to TCP connections and the incast problem [41], [42]. Referring back to Fig. 5 (solid line), with too many partitions ($k > 15$), the networking overhead outweighs the benefits of improved load balancing. To quantify the network overhead with a large number of partitions, we further increased the value of k to 100. Specifically, we placed all file partitions on the same cache server so that the total network bandwidth is kept unchanged with different partition numbers. We measured the network goodput, i.e., the transmission rate of *useful* bits excluding protocol overhead bits, with different partition numbers. Fig. 6 shows the results normalized by the goodput with one partition. When the network bandwidth is set to 1 Gbps, the goodput dropped by 20% with 20 partitions and by nearly 40% with 100 partitions. The loss of goodput is exactly the network overhead caused by reading from multiple partitions. We made similar observations when further throttling the bandwidth to 500 Mbps, where the normalized goodput gradually decreases and drops to 0.6 with 100 partitions.

Second, simple partition is susceptible to *stragglers*, as reading from many servers in parallel is bottlenecked by the slowest machine. To illustrate this problem, we manually injected stragglers into the cluster. Specifically, for each partition read, we slept the server thread with probability 0.05 and delayed the read completion by a factor randomly drawn from the distribution profiled in the Microsoft Bing cluster trace [43]. We measured the average read latency and depict the results as the dashed line in Fig. 5. As the read parallelism k increases, the latency caused by stragglers quickly dominates, leading to even longer delay. Table 3 shows the coefficient of variation in this set of experiments.

We observe that simple partition is vulnerable to straggler effects, which lead to high variances with large partition numbers.

4.3 Fixed-Size Chunking and Its Problems

Fixed-size chunking is a common practice for many distributed storage/caching systems, e.g., HDFS [18], Windows Azure Storage [17], and Alluxio [15]. With fixed-size chunking, files are split into multiple chunks of a pre-specified size and distributed across servers. However, fixed-size chunking suffers from the same problems of simple partition, as the chunk size is uniform for all files in the caching cluster. Configuring a large chunk size cannot eliminate hotspots, while a small chunk size results in too many file chunks, increasing network overhead and aggravating stragglers.

In light of these problems, we wonder: is it possible to achieve load-balanced, redundancy-free cluster caching using file splitting while still being resilient to stragglers? We give an affirmative answer in the following sections.

5 SP-CACHE: DESIGN AND ANALYSIS

In this section, we present SP-Cache, a load balancing scheme that selectively partitions hot files based on their sizes and popularities. We analyze its performance and seek an optimal operating point to minimize the average read latency without amplifying the impact of stragglers.

5.1 SP-Cache Design Overview

SP-Cache employs *selective partition* to load-balance cluster caches under skewed popularity. In a nutshell, it evenly splits a file into small partitions, where the number of partitions is *in proportion to the expected load* of that file. Specifically, for file i , let S_i be its size and P_i be its popularity. The *expected load* of file i is measured by $L_i = S_i P_i$. Let k_i be the number of partitions file i is split into. With SP-Cache, we have

$$k_i = \lceil \alpha L_i \rceil = \lceil \alpha S_i P_i \rceil, \quad (1)$$

where α is a system-wide *scale factor* applied to all files. This results in the *uniform load* across partitions, i.e., $L_i/k_i \approx \alpha^{-1}$.

SP-Cache *randomly places* k_i partitions across N servers in the cluster, where *no two partitions* are cached in the same server. Random placement improves load balancing. It ensures each server to store *approximately an equal number* of partitions. Given the uniform load across partitions, each server is expected to have the balanced load.

5.2 Benefits

SP-Cache is more efficient than simple partition (Sec. 4). It differentiates the *vital few* from the *trivial many*, in that a small number of hot, large files (*vital few*) are subject to *finer-grained* splitting than a large number of cold, small files (*trivial many*). As the former is the main source of congestion, spreading their load to more partitions mitigates the congestion on the hot spots more than doing so to the latter. Moreover, given the small population of hot files [9], [10], splitting them results in fewer partitions than splitting a large number of cold files. This, in turn, results in a reduced

number of concurrent TCP connections, alleviating the incast problem [41], [42].

Moreover, we show that SP-Cache achieves better load balancing than EC-Cache [8]. In particular, we denote by X^{SP} the total load on any particular server under SP-Cache, where X^{SP} is a random variable. Let X^{EC} be similarly defined for EC-Cache. We use the variance of X^{SP} (X^{EC}) to measure the degree of load imbalance—a higher variance implies the more severe load imbalance. The following theorem holds.

Theorem 1. Consider SP-Cache with scale factor α and EC-Cache with a (k, n) erasure code. In a cluster where the number of servers is much greater than the number of partitions of any particular file under the two schemes, we have

$$\frac{\text{Var}(X^{EC})}{\text{Var}(X^{SP})} \rightarrow \frac{\alpha \sum_i L_i^2}{k \sum_i L_i}. \quad (2)$$

Proof: Consider any particular server. Denote X_i as the load contributed by file i to this server. We have

$$X = \sum_i X_i,$$

Assuming independent partition placement across files, we have

$$\text{Var}(X) = \sum_i \text{Var}(X_i). \quad (3)$$

To facilitate the derivation of $\text{Var}(X_i)$, we define a binary random variable a_i indicating whether the request for file i is served by this server. The load X_i can be expressed as

$$X_i = a_i \frac{L_i}{k_i},$$

where k_i denotes the (non-parity) partition number of file i and $\frac{L_i}{k_i}$ calculates the partition-wise load of file i .

Under SP-Cache, each server has a probability of $\frac{k_i^{SP}}{N}$ to be selected to cache the partitions of file i . Therefore, a_i^{SP} follows Bernoulli distribution with parameter $\frac{k_i^{SP}}{N}$. We have

$$\text{Var}(X_i^{SP}) = \left(\frac{L_i}{k_i^{SP}}\right)^2 \text{Var}(a_i^{SP}) = \left(\frac{L_i}{k_i^{SP}}\right)^2 \frac{k_i^{SP}}{N} \left(1 - \frac{k_i^{SP}}{N}\right).$$

Under EC-Cache, each server has a probability of $\frac{n_i^{EC}}{N}$ to cache the partitions of file i ; each server caching the partitions has a probability of $\frac{k_i^{EC}+1}{n_i^{EC}}$ to serve the request. Therefore, a_i^{EC} follows Bernoulli distribution with parameter $\frac{n_i^{EC}}{N}$. $\frac{k_i^{EC}+1}{n_i^{EC}} = \frac{k_i^{EC}+1}{N}$. Similarly, we have

$$\text{Var}(X_i^{EC}) = \left(\frac{L_i}{k_i^{EC}}\right)^2 \text{Var}(a_i^{EC}) = \left(\frac{L_i}{k_i^{EC}}\right)^2 \frac{k_i^{EC}+1}{N} \left(1 - \frac{k_i^{EC}+1}{N}\right).$$

Suppose that the server number N is much larger than the partition number k_i . With (3), we have

$$\frac{\text{Var}(X^{EC})}{\text{Var}(X^{SP})} \approx \frac{\sum_i \left(\frac{L_i}{k_i^{EC}}\right)^2 \frac{k_i^{EC}+1}{N}}{\sum_i \left(\frac{L_i}{k_i^{SP}}\right)^2 \frac{k_i^{SP}}{N}} \approx \frac{\sum_i \frac{L_i^2}{k_i^{EC} N}}{\sum_i \frac{L_i^2}{\alpha N}} = \frac{\alpha}{k^{EC}} \frac{\sum_i L_i^2}{\sum_i L_i},$$

which completes the proof. \square

It is easy to show that under heavily skewed popularity, the variance bound (2) approaches $\frac{\alpha}{k} L_{\max}$, where L_{\max} measures the load of the *hottest* file, i.e., $L_{\max} = \max_i L_i$. As α and k are both constants, Theorem 1 states that compared to EC-Cache, SP-Cache improves load balancing by a factor of $O(L_{\max})$ in a large cluster.

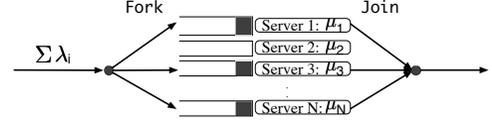


Fig. 7: The fork-join queuing model for SP-Cache.

5.3 Determining the Optimal Scale Factor

Despite the promising benefits offered by SP-Cache, it remains a challenge to judiciously determine the scale factor α . On one hand, choosing a small α results in a small number of partitions that are *insufficient* to mitigate the hot spots. On the other hand, choosing a too large α results in high I/O parallelism, adding the risk of being slowed down by stragglers.

We address this challenge with the optimal scale factor which is *large enough* to load-balance cluster caches, but also *small enough* to restrain the impact of stragglers. Specifically, we model SP-Cache as a *fork-join queue* [13], [14] and establish an *upper bound* for the mean latency as a function of scale factor α . Based on this analysis, we propose an efficient search algorithm which exponentially increases α to reduce the mean latency until the improvement becomes marginal. We settle on that α as a sweet spot, for it yields “just-enough” partitions to attain load balancing.

Model We model SP-Cache as a fork-join queue [13], [14] illustrated in Fig. 7. In particular, SP-Cache “forks” each file read to multiple reads on its partitions. Upon completion, all those partition reads “join” together to reassemble the file.

For tractable analysis, we consider Poisson arrivals of the read requests for each file. We shall verify in Sec. 7.7 that this technical assumption is not critical with real-world request arrival sequences. Let λ_i be the request rate of file i . We measure the popularity of file i as

$$P_i = \frac{\lambda_i}{\sum_j \lambda_j}. \quad (4)$$

We model each cache server as an independent $M/G/1$ queue [44] with a FIFO service discipline. We derive the mean service delay on server s , which is the partition transfer delay averaged over all reads. Specifically, let C_s be the set of files having partitions cached on server s , and B_s the available network bandwidth. For a partition of file $i \in C_s$, the transfer delay depends on its size $\frac{S_i}{k_i}$ and the network bandwidth B_s . To account for the possible network jitters, we model the transfer delay as exponentially distributed with mean $\frac{S_i}{k_i B_s}$. The chance that file i 's partition gets accessed is simply its request rate normalized by the aggregated rate, i.e., λ_i / Λ_s , where Λ_s is the aggregated request rate on server s and is given by

$$\Lambda_s = \sum_{i \in C_s} \lambda_i. \quad (5)$$

The mean service delay on server s is then computed as

$$\mu_s = \sum_{i \in C_s} \frac{\lambda_i}{\Lambda_s} \frac{S_i}{k_i B_s}. \quad (6)$$

Note that to make the analysis tractable, we assume a non-blocking network (i.e., no delay in the network) and do not model the stragglers. Our goal is to analyze the impact

of scale factor α on load balancing with respect to the mean latency.

Mean Latency We denote by $Q_{i,s}$ as the read latency file i experiences on server s , which includes both the queuing delay and the service delay (transfer delay of a partition). As the file read is bottlenecked by the *slowest* partition read, the mean read latency of file i is given by

$$\bar{T}_i = \mathbb{E}[\max_{s: C_s \ni i} Q_{i,s}]. \quad (7)$$

Summing up the mean latency over files, weighted by their popularities, we obtain the mean read latency in the system:

$$\bar{T} = \sum_i P_i \bar{T}_i. \quad (8)$$

The mean read latency critically depends on scale factor α . Intuitively, having a large α results in a large number of small partitions, which reduces both the overall queuing delay (better load balancing) and the transfer delay (small partitions).

Upper Bound Analysis Unfortunately, exactly quantifying the mean latency (8) in the fork-join system remains *intractable* due to the complex correlation between the partition placement (C_s) and the queueing dynamics [45]–[47]. Instead, we resort to establishing a tight *upper bound* to quantify the mean latency.

Prior work [45] shows that in a fork-join queue, the mean latency can be upper-bounded by solving a *convex optimization problem*. Applying this result [45, Lemma 2], we bound the mean read latency for file i as follows:

$$\begin{aligned} \bar{T}_i \leq \hat{T}_i = \min_{z \in \mathbb{R}} \left\{ z + \sum_{s: C_s \ni i} \frac{1}{2} (\mathbb{E}[Q_{i,s}] - z) \right. \\ \left. + \sum_{s: C_s \ni i} \frac{1}{2} \left[\sqrt{(\mathbb{E}[Q_{i,s}] - z)^2 + \text{Var}[Q_{i,s}]} \right] \right\}, \quad (9) \end{aligned}$$

where z is an auxiliary variable introduced to make the upper bound as tight as possible, and $\text{Var}[\cdot]$ measures the variance.

While the latency bound (9) is not closed-form, it can be efficiently computed as (9) is a convex optimization problem given the expectation and variance of latency $Q_{i,s}$. Using the Pollaczek-Khinchin transform and the moment generating function for $M/G/1$ queue [44], we have

$$\mathbb{E}[Q_{i,s}] = \frac{S_i}{k_i B_s} + \frac{\Lambda_s \Gamma_s^2}{2(1-\rho_s)}, \quad (10)$$

and

$$\text{Var}[Q_{i,s}] = \left(\frac{S_i}{k_i B_s} \right)^2 + \frac{\Lambda_s \Gamma_s^3}{3(1-\rho_s)} + \frac{\Lambda_s^2 (\Gamma_s^2)^2}{4(1-\rho_s)^2}, \quad (11)$$

where Γ_s^2 and Γ_s^3 denote the second and third moments of the service delay on server s , and ρ_s is the request intensity and is given by $\rho_s = \Lambda_s \mu_s$. Since the service delay is exponentially distributed with mean $\frac{S_i}{k_i B_s}$, we have

$$\Gamma_s^2 = \sum_{i \in C_s} \frac{\lambda_i}{\Lambda_s} \cdot 2 \left(\frac{S_i}{k_i B_s} \right)^2, \quad (12)$$

and

$$\Gamma_s^3 = \sum_{i \in C_s} \frac{\lambda_i}{\Lambda_s} \cdot 6 \left(\frac{S_i}{k_i B_s} \right)^3. \quad (13)$$

Summary: Putting it all together, given the scale factor α , we upper-bound the mean read latency as follows. We first compute the number of partitions $k_i = \lceil \alpha S_i P_i \rceil$ for each file i , based on which the expectation and variance of its read latency can be obtained, i.e., (10) and (11). Plugging them

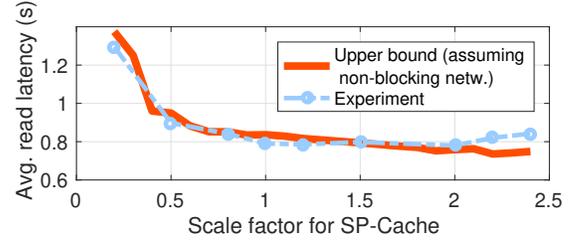


Fig. 8: Comparison of the derived upper bound and the average read latency measured in the EC2 cluster.

into (9), we solve a convex optimization problem and obtain the upper bound of the mean read latency for file i . We now upper-bound the mean latency of the system by replacing the latency of each file with its upper bound in (8).

Experimental verification: To examine how accurate the derived upper bound characterizes the mean read latency, we deployed a 31-node EC2 cluster and used it to cache 300 files (100 MB each) under skewed popularity. The detailed settings of our experiment are given in Sec. 7.1. Fig. 8 compares the derived upper bound and the mean latency measured in the cluster with various scale factors. The upper bound, though derived in a fork-join queueing model under some technical assumptions, closely tracks the average read latency measured in the EC2 cluster. Yet, as our model does not account for the networking overhead (TCP connections and incast problem) and stragglers, the measured latency occasionally goes above the theoretical upper bound.

Determining the Optimal Scale Factor We observe in Fig. 8 that as scale factor α increases, the mean latency dips quickly until an “elbow point” is reached ($\alpha = 1$), beyond which the latency plateaus for a short while and starts to rise as α turns large ($\alpha > 2$). This is by no means an accident. Intuitively, configuring a larger α results in better load balancing owing to finer-grained partitions. The price paid is the increased networking overhead and straggler impact. The gains outweigh the price before α reaches the elbow point, by which the load imbalance remains the main source of read latency. However, this is no longer the case after α passes the elbow point. The load is sufficiently balanced across servers, and the overhead of networking and stragglers becomes increasingly prominent, eventually dominating the latency.

Therefore, we should settle on the “elbow point” for the optimal scale factor. We use the derived upper bound to accurately estimate the mean latency. To locate the elbow point, we resort to an *exponential search* algorithm. Specifically, the search starts with an α such that the most heavily loaded file is split into $\frac{N}{3}$ partitions. The algorithm iteratively searches the optimal α . In each iteration, it inflates α by $1.5\times$ and examines the improvement in the derived latency bound. The search stops when the improvement drops below 1%. We shall show in Sec. 7 that using the scale factor determined by this simple algorithm, SP-Cache reduces the mean (tail) latency by up to 50% (55%) as compared to EC-Cache [8].

6 IMPLEMENTATION

We have implemented SP-Cache atop Alluxio [15], a popular in-memory storage for data-parallel clusters. In this section,

Algorithm 1 Configuration of the scale factor

```

1: procedure SP-CACHE( $\{P_i\}, \{S_i\}, \{B_s\}$ )  $\triangleright \{P_i\}$ : popularity;  $\{S_i\}$ :
   file size;  $\{B_s\}$ : network bandwidth
2:   Initialize  $\alpha^1 \leftarrow N/3/\max_i(P_i \cdot S_i)$ ,  $t \leftarrow 1$ ,  $\hat{T}^0 \leftarrow 0$   $\triangleright N$  is the
   server number;
3:    $\{C_i^t\} \leftarrow$  Random partition placement
4:   while True do
5:      $\{k_i^t\} \leftarrow$  Partition number following (1)
6:      $\{\hat{T}_i^t\} \leftarrow$  upper bound latency that solves (9)
7:      $\hat{T}^t \leftarrow$  average bound with  $\{\hat{T}_i^t\}$  following (8)
8:     if  $|\hat{T}^t - \hat{T}^{t-1}| > 0.01 \cdot \hat{T}^{t-1}$  then
9:        $t++$ 
10:       $\alpha^t = 1.5 \cdot \alpha^{t-1}$ 
11:     else
12:       break
13:   return  $\alpha^t$ 

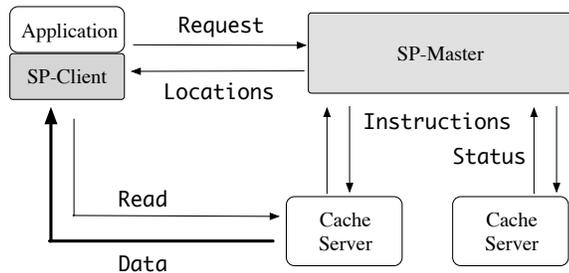
```

Algorithm 2 Parallel Repartition

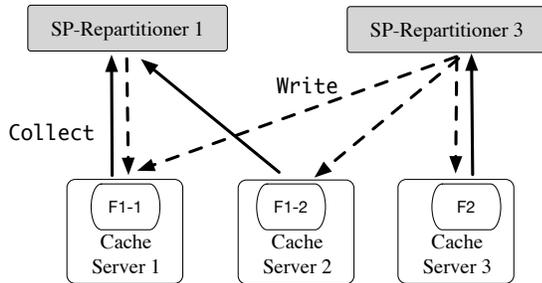
```

1: procedure PARALLEL-REPARTITION( $\{P_i\}, \{S_i\}, \{B_s\}, \{k_i^t\}$ )
2:    $\triangleright \{P_i\}$ : popularity;  $\{k_i^t\}$ : previous partition number
3:    $\alpha \leftarrow$  SP-Cache( $\{P_i\}, \{S_i\}, \{B_s\}$ )
4:    $\{k_i\} \leftarrow$  Partition number following (1)
5:    $C_s \leftarrow \emptyset$  for all server  $s$   $\triangleright$  Initialize server loads
6:   for all file  $i$  do  $\triangleright$  Keep un-repartitioned files
7:     if  $k_i \neq k_i^t$  then
8:       for all machine  $s$  storing partitions of file  $i$  do
9:          $C_s \leftarrow C_s \cup \{i\}$   $\triangleright$  Record loads
10:  for all file  $i$  do
11:    if  $k_i \neq k_i^t$  then
12:      while  $k_i > 0$  do
13:         $s \leftarrow \arg \min_{s: i \in C_s} \|C_s\|$   $\triangleright$  Server with the least load
14:         $C_s \leftarrow C_s \cup \{i\}$ 
15:         $k_i \leftarrow k_i - 1$ 
16:  return  $\{C_s\}$ 

```



(a) Read



(b) Repartition

Fig. 9: Architecture overview of SP-Cache. (a) Applications interact with SP-Client for file access. (b) SP-Repartitioners periodically repartitions files in parallel based on the instructions of SP-Master.

we describe the overall architecture of SP-Cache and justify the key design decisions made in the implementation.

6.1 Architecture Overview

Fig. 9 gives an overview of the system architecture. SP-Cache consists of three components: SP-Master, SP-Client, and SP-Repartitioner. The SP-Master implements the main logic of selective partition described in Sec. 5. It oversees many Alluxio cache servers and maintains the metadata (e.g., popularity) of files stored on those servers. The SP-Client, on the other hand, accepts the read/write requests from applications and interacts with cache servers for partition collecting and file reassembling. The SP-Repartitioners run in parallel on cache servers and re-balance the loads across servers periodically based on the instructions issued by the SP-Master (more details in Sec. 6.2).

Reads Fig. 9(a) shows the data flow for reads. Similar to Alluxio [15], an application submits read requests to an SP-Client through the provided API. Upon receiving a file request, the SP-Client contacts the SP-Master, who returns a list of cache servers containing the partitions of the requested file. The master also updates the access count for the requested file, so as to keep track of the file popularity for future *load re-balancing*. The SP-Client then communicates with the servers in the list and reads file partitions *in parallel*. Upon the completion of parallel reading, the client reassembles partitions to recover the original file and passes it to the application.

Writes SP-Cache directly writes a new file to a *randomly* selected cache server *without splitting*, given that cold files usually dominate in population (Sec. 2). For each file stored in the cluster, SP-Cache keeps track of its popularity and periodically adjusts the number of partitions based on its load: once the file turns hot, it will get repartitioned. We elaborate on how this can be done in the next subsection.

6.2 Periodic Load Balancing with Parallel repartition

Periodic repartition As file popularities may change over time, SP-Cache *periodically load-balances* cache servers by *repartitioning* the stored files. Following the recommendations in [9], SP-Cache repartitions files every 12 hours based on the access count measured in the past 24 hours. To do so, the SP-Master instructs each cache server to report its current network bandwidth (measured through *sample reads*). Based on the bandwidth and the popularity information, the master computes the optimal scale factor α using the method described in Sec. 5.3. Specifically, our implementation uses CVXPY [48] to solve Problem (9).

The effectiveness of periodic load balancing is supported by the evidence that the file popularity in production clusters is *relatively stable* in a short term (e.g., days). In fact, it has been observed in a Microsoft cluster that around 40% of the files accessed on any given day were also accessed four days before and after [9]. A similar conclusion can also be drawn from the Yahoo! cluster trace [33]: nearly 27% of the files remain hot for more than a week.

Parallel repartition SP-Cache employs two special designs to minimize the latency of load re-balancing. First, to reduce the repartition overhead, the SP-Master identifies the files

whose partition numbers remain *unchanged* and keeps them untouched in their original locations. At the same time, the SP-Master will record the load contributed by these files in each cache server such that the load distribution could be balanced in the following step.

Next, SP-Cache launches one SP-Repartitioner in *each* cache server. To facilitate parallel repartition, each SP-Repartitioner will handle a *disjoint set* of files assigned by the SP-Master. The SP-Repartitioner assembles a file and re-splits it into the specified number of partitions. To reduce the network overhead incurred by reassembling, for each file that needs to be repartitioned, the SP-Master randomly selects a SP-Repartitioner in a cache server containing partitions of that file.

Fig. 9(b) shows an example where two files $F1$ and $F2$ are repartitioned across three cache servers. Originally, file $F1$ has two partitions cached in server 1 and 2, while file $F2$ has only one partition cached in server 3. The SP-Repartitioner in server 1 (server 3) is selected to repartition file $F1$ ($F2$), as the file has a partition cached in the server, which does not need to be transferred over the network. After the file has been assembled, the SP-Repartitioner splits it and distributes the new partitions across the cache servers. In Fig. 9(b), file $F1$ is aggregated into a single partition cached in server 1; file $F2$ is repartitioned into three partitions across the three servers.

The placement of the newly generated partitions is planned in advance by the SP-Master, using a greedy placement strategy that tries to balance the load distribution as much as possible. Specifically, for a file with a new partition number k , the SP-Master chooses k distinct servers with the *smallest* loads to place its partitions. The details of the parallel repartition scheme are presented in Algorithm 2.

6.3 Partition Placement

Prevalent caching systems simply employ random data placement (e.g., the Round-Robin policy) [18], [49], [50], which increases the risk of load imbalance due to the skewed popularity. In fact, it has been acknowledged that problematic data placement is one of the root causes of load imbalance [11]. This is no longer a problem for SP-Cache. As shown in Sec. 5.1, in SP-Cache, each partition contributes approximately the same load. Therefore, random placement is sufficient to achieve good load balancing. We will verify this point in evaluations in Sec. 7.3. Moreover, we will show in Sec. 7.4 that the greedy placement strategy (Algorithm 2) proposed for parallel repartition further improves load balancing over the random placement strategy in the presence of changing popularity of files.

6.4 Implementation Overhead

Metadata SP-Cache requires only a small amount of metadata maintained in the master node. For each file i , the SP-Master stores the partition count k_i and a list of the k_i servers containing those partitions. Compared to the file metadata maintained in Alluxio, the storage overhead is negligible.

Computational Overhead Finding the optimal scale factor α appears the main computational overhead in our implementation. Nevertheless, our evaluations show that even

with 10k files, the optimal scale factor can be configured within 90 seconds (details in Sec. 7.2). As the computation is only needed every 12 hours, its overhead can be amortized and is less of a concern.

7 EVALUATIONS

In this section, we provide comprehensive evaluations on SP-Cache through EC2 deployment and trace-driven simulations. The highlights of our evaluations are summarized as follows:

- 1) The upper-bound analysis in Sec. 5 provides a reliable guidance to search for the optimal scale factor with low computational overhead (Sec. 7.2).
- 2) With 40% less memory usage, SP-Cache reduces the average read latency by 29 – 50% (40 – 70%) and tail latency by 22 – 55% (33 – 60%) over EC-Cache (selective replication) (Sec. 7.3).
- 3) In case of popularity shifts, SP-Cache is able to re-balance the load across cache servers within 3 seconds for up to 350 files; the overhead of repartition grows very slowly as the number of files increases (Sec. 7.4).
- 4) SP-Cache is resilient to stragglers, improving the read latencies by up to 40% over EC-Cache in both the mean and the tail (Sec. 7.5).
- 5) With limited cache budget, SP-Cache achieves a higher cache hit ratio than EC-Cache (Sec. 7.6).
- 6) In terms of the average write performance, SP-Cache is $1.77\times$ faster than EC-Cache (Sec. 7.8).

7.1 Methodology

Cluster Setup We have deployed SP-Cache in an Amazon EC2 cluster with 51 r3.2xlarge instances. Each node has 8 CPU cores, 61 GB memory. We measured 1 Gbps network bandwidth between instances using iPerf. We used 30 nodes as the cache servers, each with 10 GB cache space, one node as the master, and the remaining 20 nodes as clients continuously submitting read requests as independent Poisson processes.

Skewed Popularity We configured the skewed file popularity to follow a Zipf distribution [8], [51]–[53]. Unless otherwise specified, the exponent parameter of the Zipf distribution is set to 1.05 (i.e., high skewness).

Metrics We use the mean and the tail (95th percentile) read latencies as the primary performance metrics. We calculate the *improvement of latencies* as

$$\text{Latency improvement} = \frac{D - D_{\text{SP}}}{D} \times 100\%, \quad (14)$$

where D_{SP} and D denote the latencies measured under SP-Cache and the compared scheme, respectively.

In addition, we measure the degree of load imbalance by the *imbalance factor*, defined as

$$\eta = \frac{L_{\text{max}} - L_{\text{avg}}}{L_{\text{avg}}}, \quad (15)$$

where L_{max} and L_{avg} are the maximum and average load across servers. Lower values of η imply better load balancing.

Baselines We benchmark SP-Cache against three baselines.

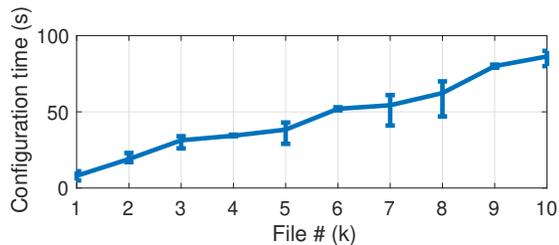


Fig. 10: Distributions of co-access distances in TPC-DS benchmark

EC-Cache: The erasure coding scheme of EC-Cache [8] determines its ability on load-balancing as well as its cache redundancy. In particular, EC-Cache claims to employ an adaptive coding strategy based on file popularities with a total memory overhead of 15%. However, the details disclosed in the paper and the implementation, kindly provided by the authors, are not sufficient for a full reconstruction of its adaptive coding scheme. Instead, we used a uniform (10, 14) erasure coding scheme for all files, which is shown to achieve the best performance in its sensitive study of coding parameters. Now that the cache redundancy is 40%, we expect EC-Cache to achieve better performance in handling the load imbalance and straggler issues.

Selective replication: For a fair comparison, we copied the top 10% popular files to 4 replicas. Therefore, assuming equal-sized files, the overall cache redundancy incurred by selective replication is $10\% \times 4 = 40\%$ —the same as that of EC-Cache.

Fixed-size chunking: Fixed-size chunking is commonly employed in many prevalent distributed storage/caching systems where files are split into multiple chunks of a constant size. Typically, the chunk size is configured as a large value, e.g., 512 MB in Alluxio. Since we test with 100 MB files in our experiments, such large chunk sizes result in one partition for each file, making no difference to load balancing. To this end, we configure fixed-size chunking with much smaller chunk sizes and compare its performance against SP-Cache (in Sec. 7.3 and Sec. 7.8).

7.2 Configuration of the Scale Factor and Partition Size

We first show that SP-Cache is able to configure the optimal scale factor α based on the derived upper bound. We ran the experiment with 300 files (100 MB each), and set the total access rate to 8 requests per second. Fig. 8 compares the derived upper bound and the average read latency measured in the experiments. Notice that in Fig. 8, we explore a larger range of α than what SP-Cache would search (Sec. 5.3) to demonstrate the tightness of the bound. We observe that the “elbow point” of the upper bound well aligns with that of the mean latency, suggesting that the upper-bound analysis can be used to accurately locate the optimal scale factor α .

Overhead The computational overhead of configuring the optimal α depends on the number of files, as it requires the latency upper bound (9) to be computed for *each* file. To quantify this overhead, we measured the runtime required to configure the optimal α in the master node with 1-10k files. Fig. 10 shows the average configuration time in 5 trials, where the error bars depict the maximum and the

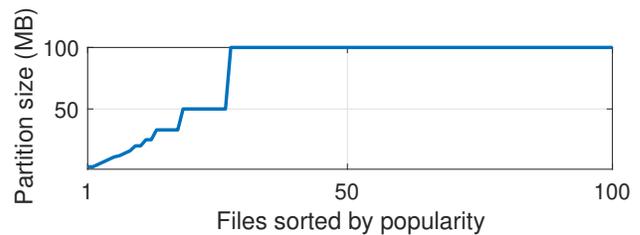


Fig. 11: The partition sizes configured by SP-Cache for files ordered by popularity (from the most popular to the least). Only the top 30% of hot files get partitioned.

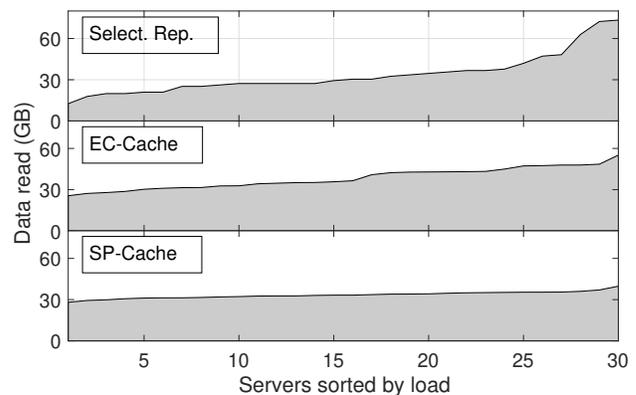


Fig. 12: Load distribution under the three load-balancing schemes. The load of a server is measured by the total amount of data reads. The client request rate is 18.

minimum. With more files, the configuration time linearly increases. Nevertheless, even with 10k files, it takes SP-Cache no more than 90 seconds to finish configuration. Since the configuration is only needed every 12 hours, its overhead is negligible.

Partition size Fig. 11 shows the optimal partition sizes SP-Cache chooses for files ordered by popularity in an experiment with 100 files (100 MB each). SP-Cache only partitions the top 30% of hot files but leaves the others untouched (no splitting). The variance in the optimal partition numbers also indicates that configuring a uniform partition number regardless of the file popularity would be highly inefficient, even with a small number of files.

7.3 Skew Resilience

We evaluated SP-Cache against the two baselines under skewed popularity in the EC2 cluster with *naturally occurred* stragglers. We cached 500 files each of size 100 MB. Note that the total cache space (300 GB) is *sufficient* to hold all 500 files and their replicas (parity partitions). The aggregated request rate from all clients is configured from 6 to 22 requests per second.

Load Balancing To study how well SP-Cache results in better load balancing than the two baselines, we measured the load of each server (i.e., the amount of data reads) under each scheme. Fig. 12 compares the load distributions under the three schemes. SP-Cache achieves the best load balancing, with imbalance factor $\eta = 0.18$. This is $2.4\times$ and $6.6\times$ better

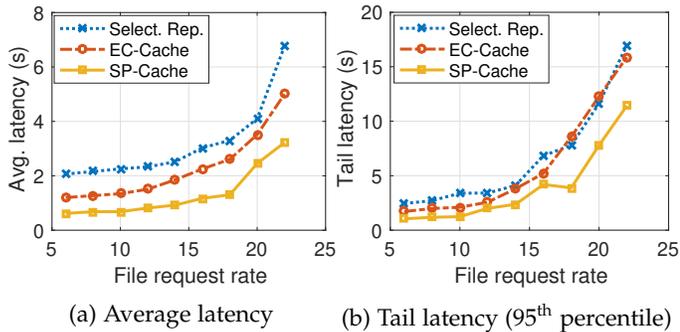


Fig. 13: Mean and tail (95th) latencies under skewed file popularity.

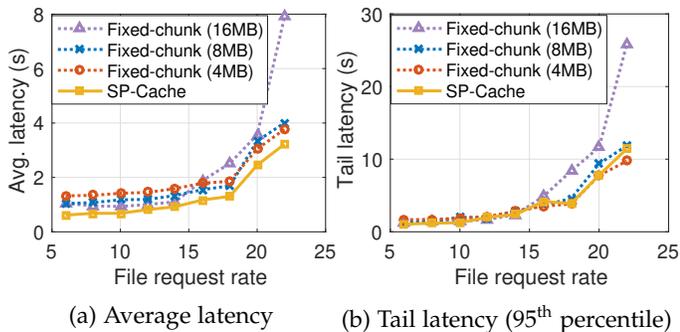


Fig. 14: Mean and tail (95th) latencies compared with fixed-size chunking.

than EC-Cache ($\eta = 0.44$) and selective replication ($\eta = 1.18$), respectively.

Read Latency Fig. 13 compares the mean and tail read latencies of the three schemes under various request rates. Owing to the improved load balancing, SP-Cache consistently outperforms the two baselines. The benefits of SP-Cache become more prominent as the request rate surges. In particular, compared to EC-Cache (selective replication), SP-Cache significantly improves the mean and tail latencies by 29-50% (40-70%) and 22-55% (33-63%), respectively.

Fixed-size chunking Fig. 14 compares SP-Cache against fixed-size chunking with chunk size of 4, 8, and 16 MB. We observe similar problems as simple partition (Sec. 4.2). On one hand, configuring small chunks results in heavy network overhead due to the increased read parallelism. We see in Fig. 14a that at low request rates (< 15), the average read latency increases as the chunk size gets smaller, e.g., up to 46% (32%) slower than SP-Cache with 4 MB (8 MB) chunks—an evidence that network overhead dominates. On the other hand, configuring large chunks, though saving the network overhead, fails to mitigate hot spots under heavy loads (request rate > 15). In fact, the mean latency with 16 MB chunks is over $2\times$ that of SP-Cache when the access rate rises to 22.

In terms of the tail latency, fixed-size chunking achieves comparable performance to SP-Cache with small chunk sizes (e.g., 4 and 8 MB), as it effectively reduces the hot spots, which are the main source of congestions in our experiments.

Compute-Optimized Cache Servers Next, we evaluate the performance of EC-Cache with cache servers of enhanced

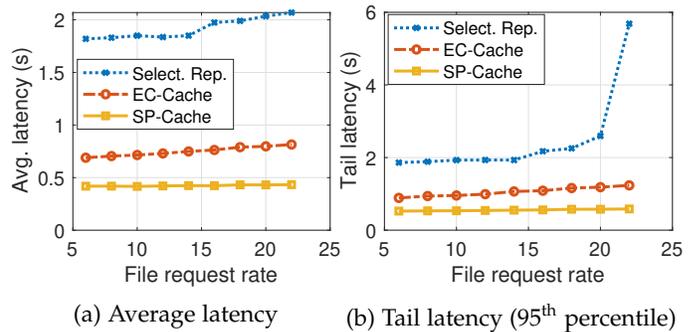


Fig. 15: Mean and tail (95th) latencies compared with c4.4xlarge (compute-optimized) instances.

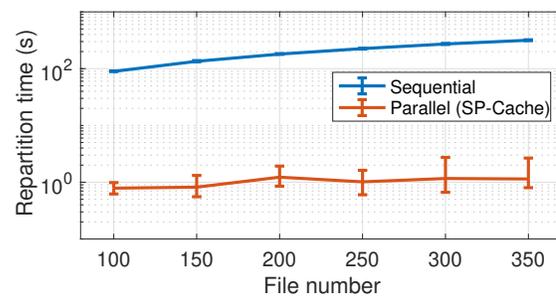


Fig. 16: Average completion time of sequential and parallel repartition schemes. Error bars depict the 95th and 5th percentiles in 10 trials.

computing power. We used c4.4xlarge instances of Amazon EC2, which are compute-optimized with Intel AVX2 extensions and the Intel Turbo Boost technology. Each c4.4xlarge instance has 16 cores and 30 GB memory. We measured 1.4 Gbps network bandwidth between c4.4xlarge instances, which is 40% higher than that in the cluster with r3.2xlarge instances.

With the boosted computing power of cache servers, we expect EC-Cache to have better performance due to the improved encoding/decoding efficiency. Fig. 15 shows the results, in which the performance gap between EC-Cache and SP-Cache remains salient. Specifically, SP-Cache outperforms EC-Cache by 39%-47% and 40%-53% in terms of the average and tail read latencies, respectively. The latency with SP-Cache remains steadily below 0.5 second on average and 0.6 second in the 95th percentile, indicating much better performance than in the cluster with r3.2xlarge instances due to higher network bandwidth. The result suggests that SP-Cache can easily handle heavy request loads with improved networking conditions. Also note that selective replication has much worse performance, leading to 3.3-3.8x and 2.5-8.7x longer latencies than SP-Cache on average and in the tail, respectively.

7.4 Resilience to Popularity Shifts

We next evaluate how SP-Cache adapts to the changing file popularities with the parallel repartition scheme. We shift the file popularity by randomly shuffling the popularity ranks of all files (under the same Zipf distribution). Notice that this presents much more drastic popularity changes than what

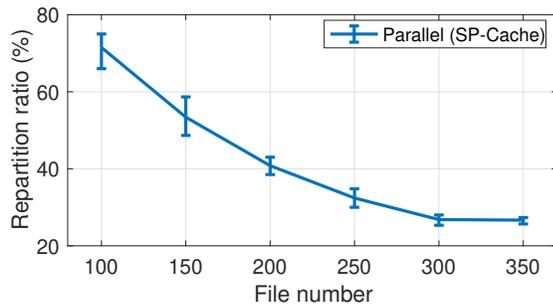


Fig. 17: Fraction of files to be repartitioned. Error bars depict the 95th and 5th percentiles in 10 trials.

has been observed in real-life clusters (Sec. 6.2) and hence imposes greater challenges of load re-balancing on SP-Cache.

We conducted experiments with files of size 50 MB each and increased the file number from 100 to 350. We measured the runtime of parallel repartition and compared it with that of a sequential scheme, where *all* files are collected and repartitioned in sequence via the SP-Master. Fig. 16 shows the results. The repartition time needed for the parallel scheme is less than 3 seconds and remains relatively stable as the number of files increases. In comparison, it takes the sequential scheme 319 seconds to finish repartition. We therefore observed a two-order-of-magnitude performance improvement achieved by SP-Cache with parallel repartition.

We attribute the substantial performance gain of SP-Cache to two main factors. First, SP-Cache speeds up the repartition process with much larger aggregated network bandwidth, as it distributes the repartition workload across many cache servers in parallel. Second, SP-Cache only needs to repartition a small fraction of files. Due to the heavy tail of the popularity distribution, most files are cold and only have a single partition. Without migrating these files, SP-Cache avoids incurring a huge volume of network traffic. Fig. 17 shows the ratio of files that need to be repartitioned after a popularity shift. We observe that the ratio decreases quickly as the number of files increases. This verifies why SP-Cache achieves a much shorter time for load re-balancing compared with the naive repartition scheme, where *all* files are collected and re-distributed sequentially.

Fig. 18 shows the load distribution under the sequential and parallel repartition schemes. Recall that the former randomly places all partitions (Sec. 5.1) while the latter employs a greedy placement (Algorithm 2) for repartitioned files. The parallel repartition scheme achieves better load balancing. By choosing the least-loaded servers to place each partition, the greedy placement algorithm can gradually improve the load balancing across servers. In addition, the overhead of the greedy searching is minimized, as only a small number of files need to be repartitioned.

7.5 Resilience to Stragglers

Redundant caching is proven resilient to stragglers [8], [9]. To evaluate how SP-Cache, which is *redundancy-free*, performs in this regard, we turn to controlled experiments with more intensive stragglers than that has been observed in the EC2 cluster. Specifically, we manually injected stragglers following the pattern profiled from a Microsoft Bing cluster

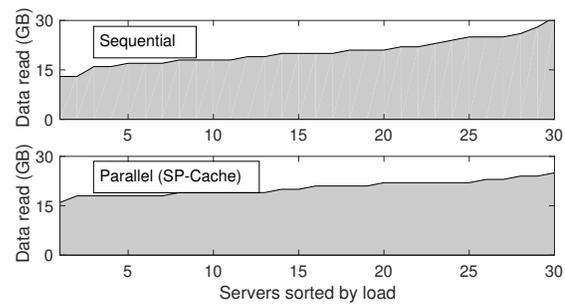


Fig. 18: Load distribution under the parallel and sequential repartition schemes. The load of a server is measured by the total amount of data reads.

trace (Sec. 4.2). We turned each cluster node to stragglers with probability 0.05 (i.e., intensive stragglers [43]).

Fig. 19 shows the results. Despite intensive stragglers, SP-Cache reduces the mean latency by up to 40% (53%) compared to EC-Cache (selective replication). Yet, the presence of stragglers results in prolonged *tail* latencies. In fact, SP-Cache exhibits slightly longer tails than the two redundant caching baselines at low request rate, as reading files from many locations adds the chance of encountering an injected straggler. As the request rate increases, most of the read requests get congested on the hot spots, and the load imbalance becomes the main source of the tail latency. Consequently, the tail latencies under the two baselines quickly ramp up. In contrast, SP-Cache effectively tames load imbalance across servers, reducing the tail by up to 41% (55%) over EC-Cache (selective replication).

7.6 Hit Ratio with Throttled Cache Budget

We stress that the benefits of SP-Cache evaluated so far were realized with 40% less memory than the two baselines. Should the same cache budget be enforced, SP-Cache would have attained even more significant benefits. To this end, we *throttled* the cluster caches and measured the *cache hit ratio* under the three load-balancing schemes. Specifically, we refer to the cluster settings in Sec. 7.3 and used the LRU (least-recently-used) policy for cache replacement. Fig. 20 compares the cache hit ratio of the three schemes with various cache budget. Owing to the redundancy-freeness, SP-Cache keeps the most files in memory and achieves the highest cache hit ratio. In comparison, selective replication falls short, as caching multiple replicas of hot files requires to evict the same number of other “not-so-hot” files out of the memory.

7.7 Trace-driven Simulation

Previous evaluations have assumed the uniform file size and Poisson arrivals of the read requests. We next remove these assumptions through trace-driven simulations with the real-world size distribution and request arrivals.

Workload We synthesized the workload based on the file size distribution and the request arrivals from two public traces. Specifically, our simulation generated 3k files. The file sizes follow the distribution in the Yahoo! traces [33] (Fig. 1); the file popularity follows a Zipf distribution with exponent 1.1. We assume that a larger file is more popular than a

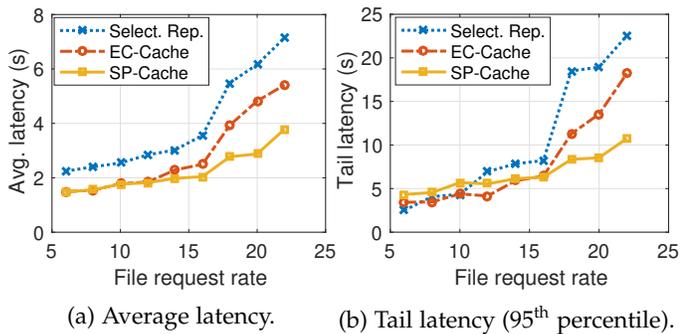


Fig. 19: Mean and tail (95th) latencies with injected stragglers.

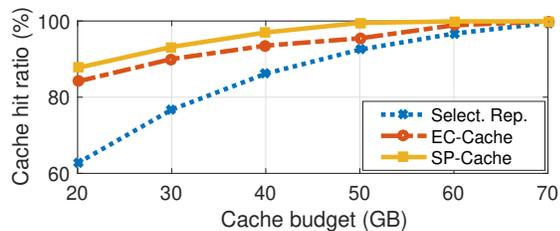


Fig. 20: Cache hit ratio with throttled cache budget.

smaller one. As the Yahoo! trace [33] provides no request arrival information, we refer to the Google cluster trace [54] which contains the submission sequence of over 660k Google cluster jobs (e.g., MapReduce and Machine Learning). Since cluster jobs usually read input at the beginning, we simply use the job submission sequence as the read request arrivals.

Settings We simulated a cluster of 30 cache servers, each with 10 GB memory and 1 Gbps network bandwidth. We manually injected stragglers into the simulated cluster as described in Sec. 7.5. We assume that a cache miss causes 3× longer read latency than a cache hit. We used a (10, 14) coding scheme in EC-Cache and set the decoding overhead as 20% (Sec. 3.2).

Results Fig. 21 shows the distributions of the read latencies under the three load-balancing schemes. SP-Cache keeps in the lead with the mean latency of 3.8 s. In comparison, the mean latencies measured for EC-Cache and selective replication are 6.0 s and 44.1 s, respectively. As hot files have large sizes in production clusters, redundant caching results in even lower cache utilization, inevitably harming its I/O performance.

7.8 Write Latency

Finally, we evaluate the write latency. We wrote files of various sizes to the EC2 cluster. We configured SP-Cache to enforce file splitting upon write based on the provided file popularity. Fig. 22 compares the write performance of the four schemes. Selective replication is the slowest, as writing multiple replicas transfers a large volume of data over the network. Using erasure code, EC-Cache writes less amount of data, but the encoding overhead drags it down. Such encoding overhead gets more significant as the file size increases, which is in accordance with the our observation in Sec. 3.2. Similarly, fixed-size chunking incurs much higher network overhead when files get larger. The reason is that

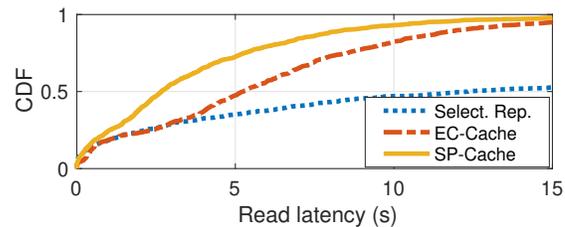


Fig. 21: Distributions of the read latencies under three load-balancing schemes in trace-driven simulations.

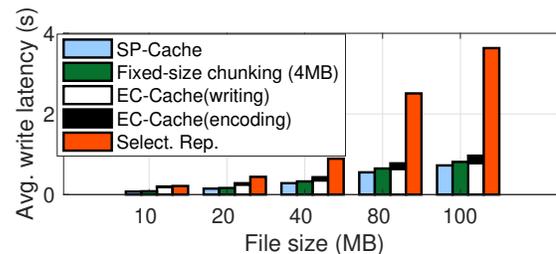


Fig. 22: Comparison of write latencies with different file sizes.

the number of chunks gradually increases and the cost of TCP connections becomes non-trivial. As shown in Fig. 6, the network goodput sharply drops with a large number of chunks.

Without suffering from these problems, SP-Cache provides the fastest writes and is on average 1.77× and 3.71× faster than EC-Cache and selective replication, respectively. Compared with fixed-size chunking using a chunk size of 4 MB, the performance advantage of SP-Cache is 13% on average. Notice that we employed sequential write in SP-Cache for fair comparison with the other three baselines. The write performance can be further improved using the parallel partition scheme of SP-Cache as described in Sec. 6.2.

8 LIMITATION AND DISCUSSION

While SP-Cache significantly outperforms the existing solutions, our current implementations have several limitations. We discuss those limitations and leave them for future explorations.

Short-term Popularity Variation With periodic load balancing, SP-Cache is unable to timely handle the short-term popularity shifts, e.g., bursts of access to certain files. To address this problem, we can enable *online and dynamic adjustment of partition granularity* in case that some files may turn hot (or cold) in a short period of time. We expect SP-Cache to respond to popularity variations much faster than EC-Cache and selective replication. To quickly adjust the partition granularity in an online fashion, SP-Cache can split and combine the existing partitions. This can be done in a *distributed manner* and incurs only a small amount of data transfer. In contrast, EC-Cache needs to collect *all* the partitions at the master node for re-encoding; selective replication incurs 1× bandwidth and storage overhead for every additional replica.

Fault Tolerance While SP-Cache manages to minimize the impact of stragglers, it does not provide fault tolerance for *non-transient* stragglers which can be *arbitrarily* slow to

the extent of a complete failure. We stress that such fault tolerance *cannot* be achieved *without cache redundancy* [8], [9]. Nevertheless, since the underlying storage system readily handles storage faults (e.g., the cross-rack replication of HDFS [18] and S3 [16]), SP-Cache can always recover the lost data from stable storages relying on the checkpointing and recomputing mechanism of Alluxio. First, Alluxio periodically persists the cached files to underlying storages (i.e., checkpointing [2]). The persisted data will then be replicated across the storage systems to ensure fault tolerance. Second, Alluxio keeps a lineage for each file, recording how this file can be recomputed from source files stored in stable storages. Once a file that has not been checkpointed gets lost, SP-Cache will recompute this file based on its lineage.

Finer-Grained Partition For structured data with clear semantics, e.g., Parquet files [55], it is unnecessary to partition or replicate the entire file uniformly if there are discrepant popularities within the file. In this case, SP-Cache can be extended to support finer-grained partition within a file by examining the popularities of different parts of the file.

9 RELATED WORK

Cluster caching has been broadly employed in data-intensive clusters as disk I/O remains the primary performance bottleneck for data analytics [4], [6], [10]. To achieve load-balanced caching, various techniques have been proposed, including data placement optimizations and replication/partition schemes.

Data placement One common approach for load balancing is to optimize the data placement scheme by designing the mapping function from files to servers. For instance, consistent hashing [50], [56] is a popular choice to implement such mappings. Unfortunately, hashing schemes may impose significant disparity on the load distribution, e.g., some heavily-loaded servers are assigned files twice as average [11]. This issue can be partly alleviated via adaptively adjusting the hash space boundaries [37], [57]. However, even with “perfect hashing” where each server holds exactly the same number of files, load balancing is not guaranteed as hashing schemes are agnostic to the skewed file popularity. Unlike these works, SP-Cache *obviates* the need for placement optimizations by eliminating the skew in the per-partition load (Sec. 5.1). The server load can then be balanced with *random placement*.

Replication Replication has been the *de facto* load balancing technique used in the disk-based object stores, including Amazon S3 [16], OpenStack Swift [31], and Windows Azure Storage [17]. Given the skewed popularity, replicating all files *uniformly* wastes the storage capacity. Selective replication [9], [58] comes as a solution. However, as popular files often have large sizes, selective replication incurs high memory overhead, and is ruled out as a practical solution for cluster caching.

File Partition EC-Cache [8] is the work most related to SP-Cache, which also takes advantage of file partition to load-balance cache servers. SP-Cache is by no means a “coding-free” version of EC-Cache. Instead, it judiciously determines the partition number of a file based on its load contribution, whereas EC-Cache simply settles on a uniform partition

scheme. To our knowledge, SP-Cache is the first work that systematically explores the benefits of selective partition.

10 CONCLUSIONS

In this paper, we have designed, analyzed, and developed SP-Cache, a load-balanced, redundancy-free cluster caching scheme for data-parallel clusters. SP-Cache selectively splits hot files into multiple partitions based on their sizes and popularities, so as to evenly spread the load of their read requests across multiple servers. We have established an upper-bound analysis to quantify the mean latency, and used it to guide the search of the optimal partition number for each file. SP-Cache effectively eliminates the hot spots while keeping the impact of stragglers to the minimum. We have implemented SP-Cache atop Alluxio. EC2 deployment and trace-driven simulations showed that SP-Cache significantly outperforms existing solutions with better load balancing in a broad range of settings. Notably, with 40% less memory footprint than EC-Cache, SP-Cache improves both the mean and the tail latencies by up to 40%, even in the presence of intensive stragglers.

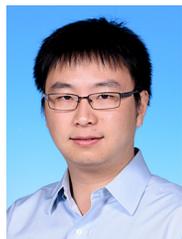
REFERENCES

- [1] Y. Yu, R. Huang, W. Wang, J. Zhang, and K. B. Letaief, “Sp-cache: Load-balanced, redundancy-free cluster caching with selective partition,” in *Int. Conf. High Performance Comput., Netw., Storage and Analysis (SC)*. IEEE/ACM, 2018.
- [2] H. Li, A. Ghodsi, M. Zaharia, S. Shenker, and I. Stoica, “Tachyon: Reliable, memory speed storage for cluster computing frameworks,” in *Proc. ACM SoCC*, 2014.
- [3] Presto, “<https://prestodb.io/>.”
- [4] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica, “Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing,” in *Proc. USENIX NSDI*, 2012.
- [5] R. Power and J. Li, “Piccolo: Building fast, distributed programs with partitioned tables,” in *Proc. USENIX OSDI*, 2010.
- [6] Memcached, “<https://memcached.org/>.”
- [7] MemSQL, “<http://www.memsql.com/>.”
- [8] K. Rashmi, M. Chowdhury, J. Kosaian, I. Stoica, and K. Ramchandran, “EC-Cache: Load-balanced, low-latency cluster caching with online erasure coding,” in *Proc. USENIX OSDI*, 2016.
- [9] G. Ananthanarayanan, S. Agarwal, S. Kandula, A. Greenberg, I. Stoica, D. Harlan, and E. Harris, “Scarlett: coping with skewed content popularity in MapReduce clusters,” in *Proc. ACM Eurosys*, 2011.
- [10] G. Ananthanarayanan, A. Ghodsi, A. Wang, D. Borthakur, S. Kandula, S. Shenker, and I. Stoica, “PACMan: coordinated memory caching for parallel jobs,” in *Proc. USENIX OSDI*, 2012.
- [11] Q. Huang, H. Gudmundsdottir, Y. Vigfusson, D. A. Freedman, K. Birman, and R. van Renesse, “Characterizing load imbalance in real-world networked caches,” in *ACM HotNets*, 2014.
- [12] Intel Storage Acceleration Library, “<https://github.com/01org/isa-l>.”
- [13] R. Nelson and A. N. Tantawi, “Approximate analysis of fork/join synchronization in parallel queues,” *IEEE Trans. Computers*, vol. 37, no. 6, pp. 739–743, 1988.
- [14] C. Kim and A. K. Agrawala, “Analysis of the fork-join queue,” *IEEE Trans. Computers*, vol. 38, no. 2, pp. 250–255, 1989.
- [15] Alluxio, “<http://www.alluxio.org/>.”
- [16] Amazon S3, “<https://aws.amazon.com/s3>.”
- [17] Windows Azure Storage, “<https://goo.gl/RqVNmB>.”
- [18] K. Shvachko, H. Kuang, S. Radia, and R. Chansler, “The Hadoop distributed file system,” in *Proc. IEEE Symp. Mass Storage Syst. and Technologies*, 2010, 2010.
- [19] Gluster File System, “<https://www.gluster.org/>.”
- [20] Amazon Elastic Compute Cloud, “<https://aws.amazon.com/ec2/>,” 2016.

- [21] A. Singh, J. Ong, A. Agarwal, G. Anderson, A. Armistead, R. Bannon, S. Boving, G. Desai, B. Felderman, P. Germano *et al.*, "Jupiter rising: A decade of clos topologies and centralized control in Google's datacenter network," in *Proc. ACM SIGCOMM*, 2015.
- [22] Huawei. NUWA. https://www.youtube.com/watch?v=OsmZBRB_OSw.
- [23] K. Asanovic and D. Patterson, "FireBox: A hardware building block for 2020 warehouse-scale computers," in *USENIX FAST*, 2014.
- [24] D. Alistarh, H. Ballani, P. Costa, A. Funnell, J. Benjamin, P. Watts, and B. Thomsen, "A high-radix, low-latency optical switch for data centers," *SIGCOMM Comput. Commun. Rev.*, vol. 45, no. 4, pp. 367–368, 2015.
- [25] C. Scott. Latency trends. <http://colin-scott.github.io/blog/2012/12/24/latency-trends/>.
- [26] IEEE P802.3ba 40 Gbps and 100 Gbps Ethernet Task Force. <http://www.ieee802.org/3/ba/>.
- [27] S. Han, N. Egi, A. Panda, S. Ratnasamy, G. Shi, and S. Shenker, "Network support for resource disaggregation in next-generation datacenters," in *ACM HotNets*, 2013.
- [28] P. X. Gao, A. Narayan, S. Karandikar, J. Carreira, S. Han, R. Agarwal, S. Ratnasamy, and S. Shenker, "Network requirements for resource disaggregation," in *Proc. USENIX OSDI*, 2016.
- [29] G. Ananthanarayanan, A. Ghodsi, S. Shenker, and I. Stoica, "Disk-locality in datacenter computing considered irrelevant," in *ACM HotOS*, 2011.
- [30] E. Jonas, S. Venkataraman, I. Stoica, and B. Recht, "Occupy the cloud: Distributed computing for the 99%," in *Proc. ACM SoCC*, 2017.
- [31] OpenStack Swift, "<https://www.swiftstack.com>."
- [32] Redis, "<http://redis.io>."
- [33] Yahoo! Webscope Dataset, "<https://goo.gl/6CZZCF>."
- [34] S. Kandula, S. Sengupta, A. Greenberg, P. Patel, and R. Chaiken, "The nature of data center traffic: measurements & analysis," in *Proc. ACM IMC*, 2009.
- [35] M. Chowdhury, S. Kandula, and I. Stoica, "Leveraging endpoint flexibility in data-intensive clusters," in *Proc. ACM SIGCOMM*, 2013.
- [36] A. Greenberg, J. R. Hamilton, N. Jain, S. Kandula, C. Kim, P. Lahiri, D. A. Maltz, P. Patel, and S. Sengupta, "VL2: a scalable and flexible data center network," in *Proc. ACM SIGCOMM*, 2009.
- [37] Y.-J. Hong and M. Thottethodi, "Understanding and mitigating the impact of load imbalance in the memory caching tier," in *Proc. ACM SoCC*, 2013.
- [38] C. Huang, H. Simitci, Y. Xu, A. Ogus, B. Calder, P. Gopalan, J. Li, S. Yekhanin *et al.*, "Erasure coding in windows azure storage." in *Proc. USENIX ATC*, 2012.
- [39] 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. VLDB Endowment*, vol. 6, no. 5, 2013, pp. 325–336.
- [40] I. S. Reed and G. Solomon, "Polynomial codes over certain finite fields," *Journal of the society for industrial and applied mathematics*, vol. 8, no. 2, pp. 300–304, 1960.
- [41] M. Alizadeh, A. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan, "Data center TCP (DCTCP)," in *Proc. ACM SIGCOMM*, 2010.
- [42] H. Wu, Z. Feng, C. Guo, and Y. Zhang, "ICTCP: Incast congestion control for tcp in data-center networks," *IEEE/ACM Trans. Netw.*, vol. 21, no. 2, pp. 345–358, 2013.
- [43] G. Ananthanarayanan, S. Kandula, A. G. Greenberg, I. Stoica, Y. Lu, B. Saha, and E. Harris, "Reining in the outliers in Map-Reduce clusters using Mantri." in *Proc. USENIX OSDI*, 2010.
- [44] B. Gnedenko and I. Kovalenko, *Introduction to Queuing Theory. Mathematical Modeling*. Birkhaeuser Boston, Boston, 1989.
- [45] Y. Xiang, T. Lan, V. Aggarwal, and Y.-F. R. Chen, "Joint latency and cost optimization for erasure-coded data center storage," *IEEE/ACM Trans. Netw.*, vol. 24, no. 4, pp. 2443–2457, 2016.
- [46] G. Joshi, Y. Liu, and E. Soljanin, "On the delay-storage trade-off in content download from coded distributed storage systems," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 5, pp. 989–997, 2014.
- [47] M. Fidler and Y. Jiang, "Non-asymptotic delay bounds for (k, l) fork-join systems and multi-stage fork-join networks," in *Proc. IEEE INFOCOM*, 2016.
- [48] CVXPY, "<http://www.cvxpy.org/>."
- [49] Alluxio Configuration, "<http://bit.ly/2jdjXTd>."
- [50] A. Lakshman and P. Malik, "Cassandra: a decentralized structured storage system," in *Proc. ACM SIGOPS Operating Systems Review*, vol. 44, no. 2, 2010, pp. 35–40.
- [51] U. Niesen and M. A. Maddah-Ali, "Coded caching with nonuniform demands," *IEEE Trans. Inform. Theory*, vol. 63, no. 2, pp. 1146–1158, 2017.
- [52] C. Roadknight, I. Marshall, and D. Vearer, "File popularity characterisation," in *Proc. ACM Sigmetrics*, 2000.
- [53] A. Clauset, C. R. Shalizi, and M. E. Newman, "Power-law distributions in empirical data," *SIAM review*, vol. 51, no. 4, pp. 661–703, 2009.
- [54] C. Reiss, A. Tumanov, G. R. Ganger, R. H. Katz, and M. A. Kozuch, "Heterogeneity and dynamicity of clouds at scale: Google trace analysis," in *Proc. ACM SoCC*, 2012, p. 7.
- [55] Apache Parquet. <https://parquet.apache.org>.
- [56] D. Karger, A. Sherman, A. Berkheimer, B. Bogstad, R. Dhanidina, K. Iwamoto, B. Kim, L. Matkins, and Y. Yerushalmi, "Web caching with consistent hashing," *Computer Networks*, vol. 31, no. 11-16, pp. 1203–1213, 1999.
- [57] J. Hwang and T. Wood, "Adaptive performance-aware distributed memory caching," in *Proc. USENIX ICAC*, 2013.
- [58] R. Nishtala, H. Fugal, S. Grimm, M. Kwiatkowski, H. Lee, H. C. Li, R. McElroy, M. Paleczny, D. Peek, P. Saab *et al.*, "Scaling Memcache at Facebook." in *Proc. USENIX NSDI*, 2013.



Yinghao Yu received the B.S. degree in Electronic Engineering from Fudan University in 2015. He is currently a Ph.D. candidate in the Department of Electronic and Computer Engineering at Hong Kong University of Science and Technology. His research interests include general resource management in big data systems, with a special focus on the performance optimization of cluster caching.



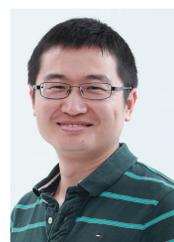
Wei Wang (S'11-M'15) received the B.Eng. (Hons.) and M.Eng. degrees from Shanghai Jiao Tong University, and the Ph.D. degree from the University of Toronto in 2015, all in the Department of Electrical and Computer Engineering. He is an Assistant Professor in the Department of Computer Science and Engineering at the Hong Kong University of Science and Technology (HKUST). He is also affiliated with HKUST Big Data Institute. His research interests cover the broad area of distributed systems, with special

emphasis on big data and machine learning systems, cloud computing, and computer networks in general.

Dr. Wang was a recipient of the prestigious Chinese Government Award for Outstanding PhD Students Abroad in 2015 and the Best Paper Runner-up Award at USENIX ICAC 2013. He was recognized as the Distinguished TPC Member of IEEE INFOCOM 2018 and 2019.



Renfei Huang received the B.Eng. degree in Computer Science and Engineering from Zhejiang University in 2018. He is currently a Ph.D. student in the Department of Computer Science and Engineering at the Hong Kong University of Science and Technology. His research interests include data visualization and distributed systems, with a special focus on the visual analytics for Industrial 4.0 and supply chain optimization.



Jun Zhang (S'06-M'10-SM'15) received the B.Eng. degree in Electronic Engineering from the University of Science and Technology of China in 2004, the M.Phil. degree in Information Engineering from the Chinese University of Hong Kong in 2006, and the Ph.D. degree in Electrical and Computer Engineering from the University of Texas at Austin in 2009. He is an Assistant Professor in the Department of Electronic and Information Engineering at the Hong Kong Polytechnic University (PolyU). His research interests

include wireless communications and networking, mobile edge computing and edge learning, distributed learning and optimization, and big data analytics.

Dr. Zhang co-authored the books *Fundamentals of LTE* (Prentice-Hall, 2010), and *Stochastic Geometry Analysis of Multi-Antenna Wireless Networks* (Springer, 2019). He is a co-recipient of the 2019 IEEE Communications Society & Information Theory Society Joint Paper Award, the 2016 Marconi Prize Paper Award in Wireless Communications (the best paper award of IEEE Transactions on Wireless Communications), and the 2014 Best Paper Award for the EURASIP Journal on Advances in Signal Processing. Two papers he co-authored received the Young Author Best Paper Award of the IEEE Signal Processing Society in 2016 and 2018, respectively. He also received the 2016 IEEE ComSoc Asia-Pacific Best Young Researcher Award. He is an Editor of IEEE Transactions on Wireless Communications and Journal of Communications and Information Networks, and an IEEE senior member.



Khaled Ben Letaief (S'85-M'86-SM'97-F'03) is an internationally recognized leader in wireless communications and networks with research interest in wireless communications, artificial intelligence, big data analytics systems, mobile edge computing, 5G systems and beyond. In these areas, he has over 600 journal and conference papers and given keynote talks as well as courses all over the world. He also has 15 patents, including 11 US patents.

Dr. Letaief is well recognized for his dedicated service to professional societies and in particular IEEE where he has served in many leadership positions, including President of IEEE Communications Society, the world's leading organization for communications professionals with headquarters in New York and members in 162 countries. He is also the founding Editor-in-Chief of IEEE Transactions on Wireless Communications and served on the editorial board of other premier journals including the IEEE Journal on Selected Areas in Communications – Wireless Series (as Editor-in-Chief).

Professor Letaief has been a dedicated teacher committed to excellence in teaching and scholarship with many teaching awards, including the Michael G. Gale Medal for Distinguished Teaching (Highest HKUST university-wide teaching award and only one recipient/year is honored for his/her contributions). He is the recipient of many other distinguished awards including the 2019 Distinguished Research Excellence Award by HKUST School of Engineering (Highest research award and only one recipient/3 years is honored for his/her contributions); 2019 IEEE Communications Society and Information Theory Society Joint Paper Award; 2018 IEEE Signal Processing Society Young Author Best Paper Award; 2017 IEEE Cognitive Networks Technical Committee Publication Award; 2016 IEEE Signal Processing Society Young Author Best Paper Award; 2016 IEEE Marconi Prize Paper Award in Wireless Communications; 2011 IEEE Wireless Communications Technical Committee Recognition Award; 2011 IEEE Communications Society Harold Sobol Award; 2010 Purdue University Outstanding Electrical and Computer Engineer Award; 2009 IEEE Marconi Prize Award in Wireless Communications; 2007 IEEE Communications Society Joseph LoCicero Publications Exemplary Award; and over 15 IEEE Best Paper Awards.

From 1990 to 1993, he was a faculty member at the University of Melbourne, Australia. Since 1993, he has been with HKUST where he has held many administrative positions, including the Head of the Electronic and Computer Engineering department. He also served as Chair Professor and HKUST Dean of Engineering. Under his leadership, the School of Engineering dazzled in international rankings (rising from # 26 in 2009 to # 14 in the world in 2015 according to QS World University Rankings). From September 2015 to March 2018, he joined HBKU as Provost to help establish a research-intensive university in Qatar in partnership with strategic partners that include Northwestern University, Carnegie Mellon University, Cornell, and Texas A&M.

Dr. Letaief received the BS degree with distinction in Electrical Engineering from Purdue University at West Lafayette, Indiana, USA, in December 1984. He received the MS and Ph.D. Degrees in Electrical Engineering from Purdue University, in Aug. 1986, and May 1990, respectively.

He is a Fellow of IEEE and a Fellow of HKIE. He is also recognized by Thomson Reuters as an ISI Highly Cited Researcher.