

# CacheSack: Theory and Experience of Google's Admission Optimization for Datacenter Flash Caches

# TZU-WEI YANG, SETH POLLEN, MUSTAFA UYSAL, ARIF MERCHANT, HOMER WOLFMEISTER, and JUNAID KHALID, Google, USA

This article describes the algorithm, implementation, and deployment experience of CacheSack, the admission algorithm for Google datacenter flash caches. CacheSack minimizes the dominant costs of Google's datacenter flash caches: disk IO and flash footprint. CacheSack partitions cache traffic into disjoint categories, analyzes the observed cache benefit of each subset, and formulates a knapsack problem to assign the optimal admission policy to each subset. Prior to this work, Google datacenter flash cache admission policies were optimized manually, with most caches using the Lazy Adaptive Replacement Cache algorithm. Production experiments showed that CacheSack significantly outperforms the prior static admission policies for a 7.7% improvement of the total cost of ownership, as well as significant improvements in disk reads (9.5% reduction) and flash wearout (17.8% reduction).

CCS Concepts: • **Information systems** → *Information retrieval*;

Additional Key Words and Phrases: Flash caches, distributed storage systems

#### **ACM Reference format:**

Tzu-Wei Yang, Seth Pollen, Mustafa Uysal, Arif Merchant, Homer Wolfmeister, and Junaid Khalid. 2023. Cache-Sack: Theory and Experience of Google's Admission Optimization for Datacenter Flash Caches. *ACM Trans. Storage* 19, 2, Article 13 (March 2023), 24 pages. https://doi.org/10.1145/3582014

https://doi.org/10.1145/3582014

# **1 INTRODUCTION**

Colossus Flash Cache (Figure 1) is the general-purpose flash cache service for Colossus [20], the successor to the Google File System [19]. Disk reads are expensive and are a major cost in datacenters: while disks are growing in storage capacity, the IO capacity (the ability to offer disk accesses per second, mainly disk reads) is not growing proportionally. As a result, to provision the IO requirements, Google needs to deploy a significant number of hard disks to serve the target IO capacity, which is costly.

The primary design goal of Colossus Flash Cache is to improve IO capacity while costing a fraction of an equivalent RAM cache or deploying more hard disks.<sup>1</sup> Colossus Flash Cache serves the read traffic of many widely used Google services including Colossus and database systems

 $^1 \rm While$  reducing read latency is also a desirable goal, it is not a design goal for Colossus Flash Cache, and beyond the scope of this article.

Authors' addresses: T.-W. Yang, M. Uysal, and A. Merchant, 1600 Amphitheatre Parkway, Mountain View, CA 94043, USA; emails: {twyang, uysal, aamerchant}@google.com; S. Pollen, H. Wolfmeister, and J. Khalid, 811 E. Washington Ave, Suite 700, Madison, WI 53703, USA; emails: {pollen, wolfmeister, junaidkhalid}@google.com.



This work is licensed under a Creative Commons Attribution International 4.0 License.

© 2023 Copyright held by the owner/author(s). 1553-3077/2023/03-ART13 https://doi.org/10.1145/3582014

ACM Transactions on Storage, Vol. 19, No. 2, Article 13. Publication date: March 2023.



Fig. 1. Colossus Flash Cache system.

such as BigQuery [37], BigTable [11], F1 [34], and Spanner [13]. Because Colossus abstracts away physical hardware complexity [20], as a Colossus service, Colossus Flash Cache only needs to focus on the aggregated cost metrics for reducing total cost of ownership, and does not need to know the actual underlying hardware.

CacheSack is the cache admission algorithm used by Colossus Flash Cache, intended to minimize the **total cost of ownership (TCO)**. Compared to a RAM cache, a flash cache usually provides a much larger cache-to-storage capacity, and so a simple algorithm such as **Least Recently Used (LRU)** may achieve a good cache hit-ratio. An idealized LRU is difficult to implement in a flash cache; we address this issue in Section 4. Flash memory has limited write endurance, so may cause premature flash wearout and increase TCO. Write amplification (Section 3.1) and flash wearout (Section 3.2), along with caching in Colossus disk servers, form a special challenge for designing a cache algorithm for Colossus Flash Cache.

#### 2 OUR CONTRIBUTIONS

CacheSack is the cache admission algorithm for Colossus Flash Cache, the successor to Lazy Adaptive Replacement Cache (LARC) [21]. CacheSack dynamically analyzes the cacheability of a workload and the given cache size, making the admission decision for the workload. CacheSack was deployed in Colossus Flash Cache in May 2021 and is now Colossus Flash Cache's default cache admission algorithm. Our contributions are summarized as follows:

- CacheSack partitions traffic into multiple categories, estimates the disk reads and cost of write of each category, and formulates a knapsack problem that finds the optimal admission policy per category to minimize the overall cost, including disk reads and bytes written to flash.<sup>2</sup>
- CacheSack effectively reduces the TCO of Colossus Flash Cache. Compared to LARC, it results in 9.5% lower disk reads, reduces bytes written to flash by 17.8%, and improves TCO by 7.7% (one week average).
- CacheSack does not require manual adjustments. That was a large engineering cost and needed when LARC was used as the admission policy.

 $<sup>^{2}</sup>$ The costs of CPU, RAM, network and power are very small relative to the cost of disk reads and bytes written to flash so they are ignored. The cost of flash storage is a fixed constant, so we omit it in the optimization.

ACM Transactions on Storage, Vol. 19, No. 2, Article 13. Publication date: March 2023.

CacheSack: Theory and Experience of Google's Admission Optimization

- CacheSack runs in real time, using a fraction of the resources of a cache index server.
- CacheSack is fully decentralized (as is Colossus Flash Cache). It requires only the information received by a single cache index server, and the failure of a single cache index server does not impact others.
- CacheSack supports major Google database systems and requires zero configuration. For other applications, users only need to provide category annotations (Section 5.1).

# 3 BACKGROUND

# 3.1 Write Amplification

Non-sequential writes to a flash drive can cause serious write amplification [29, 42], a phenomenon where one logical write causes multiple physical writes. A flash byte has to be erased before it can be rewritten. A flash block is a continuous region of bytes in a flash drive and is the smallest unit that can be erased. To erase a block, a flash drive needs to move the live bytes in the flash block somewhere else before this flash block can be erased, which causes extra writes. Write amplification reduces the IO performance and the lifetime of a flash drive; both greatly increase the total cost of ownership of a flash cache.

Sequential cache evictions (such as those caused by FIFO eviction) result in large sequential areas that can be easily erased and reused later when admitting new data. By contrast, non-sequential evictions (such as those caused by LRU eviction) result in a fragmented cache space and the flash drive has to move the interspersed live bytes somewhere else before erasing a block.

As a result, most existing eviction algorithms for RAM caches cannot be directly applied to flash caches, and write amplification is one of the most important factors to consider when designing a flash cache algorithm. Both Google [1] and Facebook [18] use FIFO-based evictions or other special purpose algorithms [36, 44] for production flash caches because of write amplification. Colossus Flash Cache reduces write amplification brought by non-sequential evictions by using *approximate* LRU (Section 4).

# 3.2 Write Endurance

Flash has limited write endurance, and thus admitting all data into Colossus Flash Cache upon write or even upon the first read would wear out the flash too soon, significantly increasing TCO. To mitigate this issue, Colossus Flash Cache previously used LARC [21] to exclude data that are accessed only once by inserting data at the second access. Figure 2 shows that more than 60% of the traffic of Colossus Flash Cache is accessed only once, and so LARC can greatly reduce bytes written to flash and avoid cache pollution.

However, excessive flash writes are still possible with LARC, and as a workaround, Colossus Flash Cache used a write rate limiter to avoid an excessively high write rate. This is, however, a blunt approach, since it does not accurately factor in the impact on overall cost, and treats all workloads similarly. It may be preferable to allow some highly cacheable workloads to burst writes at the expense of other less cacheable workloads rather than throttling all writes. CacheSack uses a more flexible and accurate approach by optimizing the total costs, including the write costs and the cost of disk reads.

## 3.3 Capturing Second-access Hits

LARC leverages the fact that a large fraction of data is accessed only once. Inserting data into cache only upon the second access avoids flash writes for data accessed once, reducing flash wear. However, the cost is that all second accesses are cache misses. Figure 2 shows that of the data accessed more than once in our workloads, 39% is accessed exactly twice, and these second accesses



Fig. 2. Fraction of bytes accessed a given number of times over a week (right truncated at 100 accesses).



Fig. 3. Miss ratios for 10 workloads at the disk server buffer cache if there is no Colossus Flash Cache (simulated).

are cache misses under LARC. This has a significant performance impact, as was also observed in Facebook's cache for social network photos [36]. Our workaround for this when using LARC was to monitor the performance loss, and to manually turn off LARC (i.e., admit all data on the first miss) for workloads that suffered a significant performance penalty. However, the manual maintenance to identify and set up special cases became more and more labor-intensive with the rapid adoption of Colossus Flash Cache in production. In our redesign, it was a requirement that the cache admission algorithm should be automatic and not require manual adjustments.

# 3.4 Colossus Buffer Cache

Colossus [20] is Google's cluster-level file system, and the next-generation of the **Google File System (GFS)** [19]. Colossus clusters scale to exabytes of storage and tens of thousands of machines. The data in Colossus is stored on "D" file servers.

In addition to Colossus Flash Cache, Colossus maintains a RAM *buffer cache* in the lower level disk D servers that buffers recent reads and writes as well as data prefetched. A cache miss in Colossus Flash Cache does not cause a disk read if the access hits in the buffer cache. Many Colossus workloads use the buffer cache extensively to improve IO performance.

In many cases, the cache hit ratio of Colossus Flash Cache is only weakly correlated with the actual disk read reduction, especially for workloads that are highly optimized for the buffer cache. Figure 3 shows simulated miss ratios of the disk server buffer caches with no Colossus Flash Cache for ten selected workloads, and they range from below 20% to over 80%. These miss ratios represent the upper bound on how far Colossus Flash Cache can improve the disk read rates. For workloads with low buffer cache miss ratios, hits in Colossus Flash Cache may simply replace buffer cache hits without improving the disk read rates. As a result, flash cache hit ratios are not a good metric to measure the efficacy of Colossus Flash Cache. In fact, our production results (Section 7.1) show that an admission policy can sometimes provide a higher hit ratio in Colossus Flash Cache but cause worse disk read rates.

# 3.5 Online and Realtime Requirements

Colossus Flash Cache is a fully decentralized system, so its cache algorithm can only use the resources of individual cache index servers, and heavyweight algorithms, such as **machine learning (ML)** models, may not be feasible. The binary of Colossus Flash Cache is updated on a weekly basis, while workloads change much more rapidly, so it is difficult for an offline-trained static model updated with the binary to adapt to workload changes. Therefore, we decided to use an online-trained model.

# 4 OVERVIEW OF COLOSSUS FLASH CACHE

Colossus Flash Cache consists of independent cache index servers. A cache index server does not directly hold cached data, but keeps an in-memory lookup table, called the index, that

tracks the locations of cached data stored in the flash drives that reside on independent storage servers.

When a Colossus Flash Cache client requests to access data stored in Colossus, the client first sends an RPC to a cache index server (see Figure 1) to determine if the requested data are already cached on a flash server (a flash hit). If so, then the cache index server sends back sufficient information for the client to access the flash copy of the data directly from the flash server. For a flash cache miss, the client contacts the disk server to read the data, while the cache index server independently decides whether to admit the data into the flash cache. If the cache index server decides to admit the data into flash, then it instructs the flash server to pull the data from the disk server directly. The extra latency of communicating with the cache index servers is negligible compared to typical remote disk read latencies, and the latencies between remote flash reads and remote disk reads are in different orders of magnitude, so Colossus Flash Cache typically reduces overall latency, although this is not an explicit service goal. The goal is reducing TCO by avoiding expensive disk reads.

The *buffer cache* of a disk server also caches recently accessed data and prefetches a small amount of data into memory for a few seconds, so that reading recently accessed data from a disk server does not necessarily cost extra disk reads. Colossus users are encouraged to design their workloads to improve IO performance by utilizing this buffer cache.

Colossus Flash Cache uses an *approximate* LRU eviction strategy to manage evictions. An idealized LRU cache would always evict the least recently used block from the cache when the cache is full. However, idealized LRU evictions cause non-sequential writes to flash, resulting in write amplification [29, 42]. To mitigate the issue of write amplification, Colossus Flash Cache uses evictions similar to Second Chance [30] to approximate LRU evictions: each cache index server manages a FIFO queue of many fixed-sized Colossus files (typically 1 GiB), each of which contains cache blocks. When evicting the file from the tail of the queue, we reinsert 28% of the most-recently used blocks into the file at the head of queue. The percentage of blocks reinserted is a tradeoff between the amount of hot blocks recycled, which improves the cache hit ratio, and the rate of reinsertion into flash, which increases write amplification. The current value (28%) is selected empirically to strike a good balance between cache performance and write amplification. This way, the write amplification factor is effectively 1.28. It is worth noting that the factor 1.28 is the software-level write amplification, which is different from the device-level write amplification. A comparison of the performance of Second Chance [30] indicates that the performance is quite close to that of LRU. Therefore, for ease of modeling, we approximate Colossus Flash Cache as an LRU cache.

Each Colossus Flash Cache server maintains a *ghost cache* [21], an in-memory lookup table that maps the key of data to the data's last access time, regardless of whether they are actually cached on flash. This is a key component of CacheSack, which relies on inter-arrival times to quickly build all the estimates described in Section 5.

Each cache index server represents a fraction of the key space, and one server's failure does not impact other cache index servers. To maintain the same reliability, CacheSack is also designed in the same decentralized manner: Each cache index server runs its own CacheSack model, using only the information received by the cache index server, and its admission decisions do not affect other cache index servers.

#### 5 CACHESACK

#### 5.1 Traffic Partitioning

CacheSack partitions potential cache blocks into many *categories*, and assigns an admission policy to each category.

The majority of Colossus Flash Cache traffic comes from Google's database systems like BigTable and Spanner where categories can be well-defined. For database traffic, CacheSack defines a category as the combination of the table name, locality group [11, 13], and type for BigTable and Spanner, and a similar combination for other databases. Since Colossus Flash Cache is also available for other Colossus users, those users can define their own categories by annotating their data. If a user does not provide a category annotation, then CacheSack will use the user name contained in the Colossus file path.

CacheSack then selects the right policy based on the pattern that category exhibits. Later, we will explain how we formulate CacheSack as a knapsack-like problem: given the cache capacity, how CacheSack chooses the items (categories) to minimize the overall cost.

#### 5.2 Admission Policies

We consider four admission policies that can be assigned to each category:

- AdmitOnWrite: Inserts a cache block at a write access or on any read cache miss.
- AdmitOnMiss: Inserts a cache block on any read cache miss, but does not insert a block at a write access. This is the conventional admission policy used in most of the cache literature.
- AdmitOnSecondMiss (LARC): Equivalent to LARC; inserts a block only after the second read access (miss), and only if the last access time is not older than the oldest last access time of the blocks in the cache, to reduce the insertion rate of cold blocks. LARC is scan resistant: Any scanned data (accessed exactly once) will not be admitted.
- NeverAdmit: Never inserts blocks.

We can sort these policies by aggressiveness: NeverAdmit < AdmitOnSecondMiss < AdmitOnMiss < AdmitOnWise.

#### 5.3 Fast Approximation to an LRU Model

To determine the best policy for the cache, the most intuitive way is to simulate all possible policycategory combinations, which is a combinatorial knapsack problem (NP-Hard). Because CacheSack currently allows up to 5,000 categories (Section 6.1) and uses four policies, there are up to  $4^{5,000}$ combinations and the knapsack problem cannot be done even with downsampled traces. Instead, we use a fast approximation for modeling an LRU cache, by introducing the *modeled cache retention time*. The cache retention time is the maximum duration that a block stays in the LRU cache without any intervening accesses to it. In practice, the cache retention time varies slowly over time. Here, we assume the *modeled* cache retention time is a constant *D* and this assumption will make all our estimates just approximations.

We use AdmitOnMiss as an example. For a given block, when a read access arrives, we can compute d, the time since last access (which is  $\infty$  if the current access is the first read). We can classify the inter-arrival times by using D (Figure 4):

- *d* ≤ *D*: An access arrives before the block leaves the cache, and therefore the access generates a cache hit and moves the block to the head of the queue.
- *d* > *D*: An access arrives after the block leaves the cache, and therefore it is a cache miss, which causes a write to the cache.

In other words, we approximate the LRU cache by a cache that has the TTL value D and resets the TTL counter of a block when receiving an access to the block. The theoretical aspect of the TTL approximation was also studied in literature. Fagin [17] showed the TTL approximation is asymptotically exact for independent and identically distributed requests, and [23] proved that given the assumption that data accesses are stationary and ergodic, the TTL approximation will



Fig. 4. LRU evictions are approximated by TTL evictions with the modeled retention time D while the TTL counter of a cache block is reset whenever the cache block is accessed. If the access interarrival time d is less than or equal to D, then this access is a cache hit, and we move the cache block to the head of the queue (the TTL counter is reset). If the access interarrival time d is greater than D, then this access is a cache miss, and we insert the cache block to the head of the queue (the TTL counter is also reset).

converge to an LRU cache as the cache size goes to infinity. The accuracy of the TTL approximation in production is analyzed in Section 7.1.

A cache miss in Colossus Flash Cache will cause a disk read if it is also a miss in the Colossus buffer cache. Each cache index server maintains a *buffer cache simulator*, and when d > D, we run the simulator and see whether it is a miss.

This way, when a new access arrives, we are able to update the disk reads, cache usage, and bytes written to flash cache caused by admitting the block using AdmitOnMiss. We can also compute the same quantities for other policies: AdmitOnSecondMiss, AdmitOnWrite, and NeverAdmit. The detailed estimation is described in Section 9.2.

A nice property of this approximation is that the estimates for a block are not affected by other blocks or policies, as long as the modeled cache retention time is given. Therefore, the disk reads, cache usage, and written bytes caused by admitting a category are just the sums of the corresponding block-level quantities.

# 5.4 Knapsack Problem

Once we have the estimates for disk reads, cache usage, and bytes written to flash cache for each policy-category pair, we have a knapsack problem: find the optimal policy per category to minimize the overall cost (disk reads, flash storage, and written bytes) while fitting within the cache. We do not specify the relative cost of disk reads, bytes written to flash, and flash storage, because they are confidential.

We further allow *fractional* policies: CacheSack can apply a policy to a fraction of a category. For example, CacheSack may decide it is optimal to apply AdmitOnMiss, AdmitOnSecondMiss, AdmitOnWrite, and NeverAdmit to 30%, 20%, 10%, and 40% of blocks in a category, respectively. Then the problem becomes a *fractional* knapsack problem [14] that finds the optimal policy fractions per category to minimize the overall cost. The advantage of considering a fractional knapsack is that it can be solved efficiently by a greedy algorithm, as opposed to a combinatorial knapsack

that is NP-Hard. Our problem is slightly different from the original fractional knapsack in [14], because we need to decide four fractions (seven fractions with additional spatial prefetch policies in Section 8) per category instead of two. Section 9.4 explains the details of how we solve our problem by a greedy algorithm after applying Andrew's monotone chain convex hull algorithm [2]. We note that if an LRU cache is perfectly modeled by the TTL approximation, the resulting cache retention time of the LRU cache is exactly D after applying the optimal policy fractions per category.

# 5.5 Optimization over Modeled Cache Retention Times

The knapsack problem in Section 5.4 is to find the optimal policy fractions for a *given* modeled cache retention time D, which cannot be known in advance. Thus, we need to solve the same knapsack problem for *all* possible D. To do this in production, we can have a set of predefined modeled cache retention times:  $0 < D_1 < D_2 < \cdots < D_m = \mathcal{D}$ , where  $\mathcal{D}$  is a suitable upper bound, and solve *m* different knapsack problems. Thanks to the greedy algorithm, we can still solve many knapsack problems (currently 127, Section 6.2) quickly.

# **6** CACHESACK IN PRODUCTION

CacheSack is now deployed in production as the default cache admission algorithm for Colossus Flash Cache. This section explains the engineering efforts needed to do so.

# 6.1 Category Assignment

The number of categories encountered in production cannot be known in advance, so we balance the need for accuracy and space by hashing a category to one of 5,000 buckets. Categories assigned to the same bucket are treated as combined in the optimization. The number of hash buckets is a trade-off between memory usage and hash collisions. The typical number of categories per server is less than 100 and our experiments showed that with 5,000 buckets, 95% of the clients see a hash collision rate lower than 1% and the worst collision rate is less than 5%. Further, cache collisions are not persistent, since each cache index server uses a different hash key and changes it periodically to break possible spatial and temporal correlations.

A bucket without sufficient training data might not provide meaningful metrics. If a bucket contributes to less than 0.1% of total lookups, then it will be aggregated to a single *catch-all* bucket before solving the knapsack problem.

# 6.2 Modeled Cache Retention Times

Currently, CacheSack uses 127 predefined cache retention times: 15 min,  $1.06 \times 15$  min,  $1.06^2 \times 15$  min,  $\ldots$ ,  $1.06^{126} \times 15$  min  $\approx 16$  days; the 128th value is reserved for positive infinity.

These retention times are decided as follows. We first determine the working range of retention times. A retention time lower than 15 min means we evict and insert cache blocks in an extremely aggressive way, which would cause serious flash wearout. By policy, any cache block is forced to leave the cache if it stays more than 15 days. Hence, we set the modeled working range of retention times as 15 min to 15 days. We then decide the number of retention times to model. We tried 127 (6% geometric increase) and 255 (3% geometric increase) retention times, and our experiments showed that 127 retention times gave similar results while reducing RAM usage by half.

# 6.3 Ghost Cache

Since LARC was Colossus Flash Cache's previous admission control, a *ghost cache* was implemented in cache index servers. It is an in-memory lookup table that maps a data's key to the data's last read access time, and LARC uses the information to determine whether to admit the data on miss. CacheSack uses the same ghost cache to obtain inter-arrival times. In addition, to build the metric estimate for AdmitOnWrite, we expanded the ghost cache so that we know whether the last access is a write access. To build the metric estimate for AdmitOnSecondMiss, we use the ghost cache to record the most two recent access times.

Because the ghost cache is the ground truth for CacheSack, the ghost cache must contain sufficient history. The optimal solution of CacheSack will not be affected as long as the ghost cache TTL, the time since the oldest last access time of the blocks in the ghost cache, is greater than the optimal modeled cache retention time. As a rule of thumb, we provision the size of the ghost cache so that its TTL is at least twice the solved optimal modeled retention time (typically about 4 h).

# 6.4 Buffer Cache Simulators

A cache miss in Colossus Flash Cache causes a disk read only if it is also a miss in the buffer cache. CacheSack simulates the buffer cache to determine whether the current miss in Colossus Flash Cache is also likely a miss in the buffer cache. In fact, we need *many* simulators: one for each pair of policy-retention time so there are 382 simulators ( $3 \times 127 + 1$ , the retention time does not affect NeverAdmit). Running the simulators is the most computationally intensive component in the CacheSack model. Fortunately, the buffer cache simulator is simple enough and only requires the access history in the past few seconds so it only moderately increases CPU load on the low-QPS servers (5% CPU usage).

# 6.5 Model Training

We use a simple scheme to train the CacheSack model: the model is reset every 5 min and is trained based on the lookups in this 5-min period. We note that a lookup contains the access times of the most recent two accesses and therefore the lookups in a 5-min period may contain the information of many hours.

The selection of the training duration is a trade-off. Using a larger training duration means the model can be improved by more training data and longer time horizon, while the model can react more quickly to changes in the workload with a shorter duration. We tested several training durations and found that 5-min one gave the best disk read reduction, although we did not find significant differences among all candidates.

# 6.6 Summary of CacheSack in Production

In summary, CacheSack uses a TTL model (Section 5.3) to quickly simulate the flash cache hits or misses based on block inter-arrival times, and uses the Colossus buffer cache simulator (Section 6.4) to simulate whether a flash cache miss incurs a disk read. The block inter-arrival times are calculated from the lookups sent by the clients to Colossus Flash Cache, and are recorded in the ghost cache (Section 6.3) of the cache index server. CacheSack solves the knapsack problem (Section 5.4) to find the optimal policy per category to minimize the TCO. The knapsack problem is solved every 5 min by using the trace that are lookups received by the cache index server within this 5-min period (Section 6.5).

## 6.7 Lessons Learned

Automatic Cache Optimization Incentivized User Adoption. In deciding whether to use Colossus Flash Cache, users weigh both the likely TCO improvement and the engineering effort required to configure and maintain it. In the past, users had to manually choose the admission policy (using AdmitOnMiss or AdmitOnSecondMiss) based on knowledge of their workload or by running A/B experiments with the assistance of the Colossus Flash Cache team. For heavy users like Spanner, Colossus Flash Cache had to provide heuristic, hand-tuned admission policies to improve cache

performance. Such human tuning and maintenance usually requires effort from both the users and the Colossus Flash Cache team, which can discourage the adoption of Colossus Flash Cache if the expected hardware resource saving does not justify the extra engineering cost.

We found that CacheSack greatly incentivized users to adopt Colossus Flash Cache. The automatic cache provisioning brought by CacheSack requires almost no configuration and maintenance so that it can be set and forgotten. We found that new users were more willing to use Colossus Flash Cache once they knew it would automatically adjust the cache policy based on their workloads.

Some of Colossus Flash Cache's existing users have independently verified that CacheSack applied appropriate admission policies to their workloads, based on the knowledge of their workloads and reporting provided by Colossus Flash Cache. One user experimentally overrode CacheSack with manually optimized policies and found that CacheSack worked as well as manual policy tuning. After CacheSack became the default admission policy in Colossus Flash Cache, we were able to retire the hand-tuned optimization for Spanner, and our existing users did not need to manually adjust the policy anymore.

*Experiment Infrastructure Accelerated Feature Development.* The development of CacheSack was significantly helped by the experiment infrastructure of Colossus Flash Cache. The experiment infrastructure allows developers to test new features by using 10% of the cache index servers, and because cache index servers are independent and isolated, any experiment can only cause minor service degradation in the worst case. Before the full deployment, we ran CacheSack as an experiment for a few months and most of the issues were identified and corrected during the experimental phase. In fact, CacheSack has caused no binary rollbacks since the full deployment.

In addition, because each server represents a fraction of the key space, which is permuted randomly, each server is statistically indistinguishable. We can have simultaneous comparisons between CacheSack and the control group to see whether CacheSack works as expected and identify any issues. The experiment infrastructure is extensively used by the developers of Colossus Flash Cache for new features, and the impact of a new feature can be accurately measured before the full deployment.

Model Introspectability and Maintainability Played Important Roles. We found that the model introspectability played an important role for the adoption of the new cache algorithm. Because any cache algorithm of Colossus Flash Cache will be operated and maintained by developers and **site reliability engineers (SREs)** after the initial deployment, one requirement of deploying a new cache algorithm is that the model behavior can be fully understood and monitored by the developers and SREs. CacheSack satisfies this requirement as it only assumes that the TTL approximation (Section 5.3) is sufficiently close to the eviction of Colossus Flash Cache, and all model behaviors can be derived from this assumption. Another advantage of a highly introspectable model is that the developers (besides the original designers) of Colossus Flash Cache can easily ensure thorough test coverage, validate software releases, and extend the original functionality of CacheSack without assistance from the original designers. After the deployment of the original CacheSack, it became the foundation of further optimizations for Colossus Flash Cache.

It is also worth mentioning that CacheSack is simple enough to be implemented by limited extensions to the original codebase of Colossus Flash Cache. In particular, the optimization was implemented as a simple greedy algorithm instead of using a generic linear program solver library. This did cost extra time for development, but we decided to do so, because it allowed us to minimize the computational overhead and increase system reliability by reducing external dependencies. More importantly, anyone familiar with the ecosystem of Colossus Flash Cache can easily maintain CacheSack or develop new features based on it. The implementation of CacheSack can evolve continuously with Colossus Flash Cache, reducing maintenance burden. Since the



Fig. 5. CacheSack disk read rate prediction errors relative to the actual value in production (CDF).



Fig. 6. Policy distribution suggested by CacheSack in selected datacenters, demonstrating a variety of workload responses.

completion of the initial deployment, involvement from the original designers has not been required for maintenance and new feature developments.

# 7 EVALUATION

#### 7.1 Production Evaluation

*Model Accuracy.* There are two LRU approximations in Colossus Flash Cache: Colossus Flash Cache uses Second-Chance-like approach to approximate LRU evictions (Section 4), and CacheSack models an LRU cache as a TTL approximation (Section 5.3). Therefore, it is important to verify that the CacheSack model is a good enough approximation to the actual Colossus Flash Cache. We examined the accuracy of CacheSack as follows. For each client, the solution to the knapsack problem in Section 5.4 gives the predicted disk reads when using the optimal admission policies. Then Colossus Flash Cache applies the optimal policies in production. We compared the predicted disk reads with the actual disk reads to see how well they match. Figure 5 shows the prediction errors of CacheSack relative to the actual values obtained from the disk servers; 51% of the relative errors are within 10% and 82% of the relative errors are within 20%.

*Policy Distribution.* Figure 6 shows the policy distributions suggested by CacheSack in the selected datacenters of various workloads. We can see that each datacenter has a different workload pattern and CacheSack adaptively decides suitable admission policies based on workloads and cache sizes. Although it would be possible for manual selection of static policies to match each datacenter workload, CacheSack is able to reduce the human toil, response delay, and operational complexity required to maintain these assignments.

Production Experiments. By using the experiment infrastructure of Colossus Flash Cache, we can compare the performance of different cache algorithms in production. Because each cache index server represents a fraction of the key space, the pattern of workload each cache index server receives is statistically indistinguishable. We let 10% of the cache index servers run static AdmitOnMiss and another 10% of the cache index servers run AdmitOnSecondMiss so that we can compare CacheSack, static AdmitOnMiss and static AdmitOnSecondMiss simultaneously in production.

From Figures 7 and 8, we see that compared to AdmitOnSecondMiss, CacheSack results in fewer disk reads (6% of one week average) and reduces 26% (one week average) written bytes to flash, and Figure 9 shows that CacheSack effectively reduces TCO in production: the cost of disk reads, flash cache writes and flash storage of CacheSack is 93% of AdmitOnSecondMiss and is 78% of AdmitOnMiss (one week average).

Figure 10 shows that CacheSack has a higher hit ratio than AdmitOnSecondMiss but lower than AdmitOnMiss. Nevertheless, AdmitOnMiss is not the best choice. Figure 7 shows that AdmitOnMiss has the worst disk read reduction even though it has the highest hit ratio. Because of the lower-level



Fig. 7. Disk reads of different admission policies in production, divided by the average value for AdmitOnSecondMiss.



Fig. 9. Total cost (a function of disk reads, flash storage and written bytes) of different admission policies in production, divided by the average value for AdmitOnSecondMiss.



Fig. 8. Written bytes of different admission policies in production, divided by the average value for AdmitOnSecondMiss.



Fig. 10. Hit ratios in Colossus Flash Cache of different admission policies in production.

buffer cache, a higher hit ratio in the flash cache does not necessarily imply fewer disk reads: many major Colossus users optimize their workloads by accessing the same data many times within the first few seconds so that only the first access causes an actual disk read. In this case, AdmitOnMiss generates many flash hits that do not reduce disk reads at all. AdmitOnSecondMiss resolves this issue by avoiding a cache insertion if the most recent access time is too recent to expect that the data has left the buffer cache.

# 7.2 Evaluation by Simulations

In addition to production experiments, we also used the *Colossus Flash Cache simulator* to test the performance of CacheSack in a variety of configurations and contexts, such as cache size and optimization iteration period. The Colossus Flash Cache simulator is used for multiple purposes including performance-regression testing by Colossus Flash Cache developers and for datacenter resource planning by Colossus Flash Cache clients. The Colossus Flash Cache simulator uses the same production code as Colossus Flash Cache, and we use production traces (sampled lookups received by cache index servers) as the input of the simulator.

We first compare the performance of CacheSack, to the static admission policies AdmitOnMiss, AdmitOnSecondMiss and AdmitOnWrite for various cache sizes. We use here a two-day trace from one large (order of million QPS) production cache as a representative. This trace reflects a uniform sample of the data accesses from a large collection of internal production workloads.

Impact of Cache Size on Performance. When the cache size is small, AdmitOnSecondMiss has a better performance than AdmitOnMiss or AdmitOnWrite, because single-use keys are excluded. However, AdmitOnMiss and AdmitOnWrite will outperform AdmitOnSecondMiss for a large cache, because second accesses will hit in the flash cache.



Fig. 11. Hit ratios in Colossus Flash Cache of different admission policies in simulation. Above: Original hit ratios. Below: Values relative to AdmitOnSecondMiss.



Fig. 12. Disk reads of different admission policies in simulation. Above: Constant scaling by dividing the values by the average value for AdmitOnSecondMiss. Below: Values relative to AdmitOnSecondMiss.

CacheSack learns to use a more conservative policy for a small cache and a more aggressive policy for a large cache. Figures 11 and 12 show that CacheSack can provide a good performance for the entire range of flash cache sizes.

It is also interesting to see the amount of written bytes caused by different admission policies in Figure 13. For AdmitOnMiss, AdmitOnWrite and AdmitOnSecondMiss with excessively small cache, blocks are frequently evicted from and reinserted into the cache, resulting in a very large amount of written bytes, especially for AdmitOnMiss and AdmitOnWrite. CacheSack, however, takes into account the cost of written bytes, and therefore only admits the most valuable part of the workload into the cache.

We can also view the total cost (a confidential function of disk reads, flash storage, and written bytes) as a function of cache size. When the cache size is small, disk reads and writes to flash are the largest contributions to cost, while flash storage is the largest cost component for larger cache sizes. Therefore, the total cost is a U-shape curve, and we are able to find the optimal cache size that minimizes the total cost. Figure 14 shows that CacheSack gives the lowest total cost for all cache sizes. CacheSack avoids the trade-off and provides robust good behavior over the range of cache sizes.

*Optimization Frequency.* We evaluated the system performance on the choice of different optimization frequencies. Here, we test different lengths of training duration from 1 min to 8 h, which span a majority of the observed time variation of workloads. Figure 15 shows that the training duration does not significantly impact the performance and all the cost metrics are similar. Because this method is insensitive to this parameter, customized or automated tuning was deemed unneeded, and the entire deployment currently uses a single value.

### 8 SPATIAL PREFETCH

## 8.1 Prefetch Policies

In addition to caching requested blocks, Colossus Flash Cache implemented *spatial prefetch* that can admit blocks that have not yet been requested but will be likely requested shortly, so that



Fig. 13. Written bytes of different admission policies in simulation. Above: Constant scaling by dividing the values by the average value for AdmitOnSecondMiss. Below: Values relative to AdmitOnSecondMiss.

ŝ

Scaled flash cache size

4

5

ż



Fig. 14. Total cost (a function of disk reads, flash storage and written bytes) of different admission policies in simulation. Above: Constant scaling by dividing the values by the average value for AdmitOnSecondMiss. Below: Values relative to AdmitOnSecondMiss.

storage users can save the seek time for these blocks. For example, to read a contiguous 1 MiB region, a hard disk typically spends approximately 10 ms for the seek and another 10 ms for the read. Therefore, if there will likely be several requests accessing the different parts of this 1 MiB region, then we can prefetch and cache the entire 1 MiB region on the first access so that storage users do not need extra seek time (which is as expensive as the read time) for the other accesses. While Colossus already uses Colossus buffer cache (Section 3.4) to prefetch and cache data in the RAM cache for a time scale of seconds, Colossus Flash Cache now can prefetch and cache data in the flash cache for minutes or hours.

Colossus Flash Cache utilizes the fact that requested blocks are of fixed-size and aligned: block offsets are multiples of the block size, and therefore blocks are non-overlapping. Colossus Flash Cache also defines a prefetch region as an aligned region (the offset is a multiple of the prefetch size), consisting of a fixed number of blocks. This makes it easy to determine the set of blocks to prefetch and whether the blocks have been cached. The alignment of prefetch regions makes sure that the admission of one prefetch region will not affect other prefetch regions, because prefetch regions are not overlapping. CacheSack provides prefetch policies analogous to AdmitOnMiss, AdmitOnSecondMiss, and AdmitOnWrite:

- PrefetchOnMiss: Similar to AdmitOnMiss, inserts the entire prefetch region on the first read access to this prefetch region.
- PrefetchOnSecondMiss: Similar to AdmitOnSecondMiss, inserts the entire prefetch region only after the second read access to the same prefetch region, and only if the last access time is not older than the oldest last access time of the blocks in the cache. Since a prefetch region is the set of aligned blocks, its access history is the combination of the access histories of the individual blocks, which can be obtained from the ghost cache.
- PrefetchOnWrite: Unlike PrefetchOnMiss and PrefetchOnSecondMiss, we admit only the requested block to the cache on the write access, and admit the entire prefetch region on any read access to this prefetch region. In other words, PrefetchOnWrite is effectively

ACM Transactions on Storage, Vol. 19, No. 2, Article 13. Publication date: March 2023.



Fig. 15. Hit ratios, disk reads, written bytes, and total cost (a function of disk reads, flash storage and written bytes) of Colossus Flash Cache with different training durations. Disk reads, written bytes, and total cost are divided by the average value for the 5-min training duration.

AdmitOnWrite and PrefetchOnMiss applied to the same category. PrefetchOnWrite does not insert the entire prefetch region on the write access, because the other blocks in the same prefetch region might not have been written to the hard disk and cannot be cached.

Since prefetch regions are non-overlapping, from CacheSack's point of view, prefetch regions are independent blocks with larger sizes, which makes the metric estimation for the prefetch policies essentially the same as their non-prefetch counterparts; the CacheSack metric estimation would be much more difficult if we were to allow overlapped prefetch regions. CacheSack relies on the last two access times to build the metric estimation (Section 9.2), and the last two access times to a prefetch region can be calculated from the combined last two access times of the blocks contained in this prefetch region, which are recorded in the ghost cache. With the prefetch feature, CacheSack evaluates the cache benefit of all policies (AOM, AOW, AOSM, POM, POW, POSM, NA) applied to each category and solves the knapsack problem to decide the optimal policy per category.

# 8.2 Production Evaluation

As in Section 7.1, we can compare CacheSack with the prefetch feature to the original, non-prefetch CacheSack in production, by running the experiment in a fraction of index servers.

Although not by a large margin, Figure 16 shows that CacheSack with prefetch consistently uses fewer disk reads than CacheSack without prefetch does. Figure 17 shows that CacheSack with the prefetch feature reduces disk reads by 3.74% (one week average), compared to the non-prefetch version of CacheSack. However, Figure 18 shows that CacheSack with prefetch also inserts data



Fig. 16. Disk reads of CacheSack with prefetch and CacheSack without prefetch in production, divided by the average value for CacheSack without prefetch.



Fig. 18. Written bytes of CacheSack with prefetch and CacheSack without prefetch in production, divided by the average value for CacheSack without prefetch.



Fig. 20. Total cost (a function of disk reads, flash storage and written bytes) of CacheSack with prefetch and CacheSack without prefetch in production, divided by the average value for CacheSack without prefetch.

to flash more aggressively, resulting in 11.04% more written bytes (one week average, Figure 19). Compared to LARC, CacheSack with prefetch reduces disk reads by 9.5% (one week average) and reduces bytes written to flash by 17.8% (one week average).

Although the spatial prefetch casues more written bytes, the net saving brought by it is still positive: Figures 20 and 21 show that CacheSack with prefetch reduces the total cost (which factors in disk read, bytes written to flash, and flash storage) by 1.33% (one week average). Taking into account the 6.5% TCO reduction brought by CacheSack without prefetch compared to LARC, CacheSack with prefetch reduces TCO by 7.7% compared to LARC.



Fig. 17. Disk reads reduction by CacheSack with prefetch, relative to CacheSack without prefetch.



Fig. 19. Extra written bytes caused by Cache-Sack with prefetch, relative to CacheSack without prefetch.



Fig. 21. Total cost reduction by CacheSack with prefetch, relative to CacheSack without prefetch.

# 9 MATHEMATICAL MODEL OF CACHESACK

#### 9.1 Model Assumption

CacheSack models the cache as using LRU evictions. Colossus Flash Cache considers data for caching to be immutable after being written, i.e., the first access is a write, and subsequent ones are reads; mutability is handled by higher layers in the system. The CacheSack model does not need the immutability assumption, but we keep it to align with the actual system; the model can be easily modified for the mutable case.

# 9.2 Metric Estimation of an LRU Cache

We begin with AdmitOnMiss. For a given block b, let  $t_1, t_2, t_3, ..., t_n$  be the read access times, and  $t_0 = -\infty$  for convenience. Therefore, the inter-arrival times are  $d_i = t_i - t_{i-1}$  and  $d_1 = t_1 - t_0 = \infty$ . Assume that D is the modeled cache retention time; that is, D is the maximum duration that a block stays in the LRU cache without any intervening accesses. We can classify the inter-arrival times as follows:

- *d<sub>i</sub>* ≤ *D*: A cache hit because the access arrives before the block leaves the cache. The block is moved to the head of the queue because of the LRU eviction.
- $d_i > D$ : A cache miss because the access arrives after the block leaves the cache. AdmitOnMiss inserts the block into the cache on miss, causing a write to the cache.
- For a flash cache miss, we update the *buffer cache simulator* to see whether it is also a cache miss in the buffer cache. If so, then the access is a disk read.

We can then write disk reads  $S_b^{AOM}(D)$ , cache byte-time usage<sup>3</sup>  $U_b^{AOM}(D)$ , and bytes written to cache  $W_b^{AOM}(D)$  as functions of D:

$$B_{b}^{\text{AOM}}(D, i) = \begin{cases} 1, & \text{Buffer Cache Hit at } t_{i}, \text{ using AOM} \\ 0, & \text{Buffer Cache Miss at } t_{i}, \text{ using AOM} \end{cases},$$
$$S_{b}^{\text{AOM}}(D) = |\{i : d_{i} > D, B_{b}^{\text{AOM}}(D, i) = 0\}|,$$
$$U_{b}^{\text{AOM}}(D) = \text{Size}(b) \times \sum_{i} \min(d_{i}, D),$$
$$W_{b}^{\text{AOM}}(D) = \text{Size}(b) \times |\{i : d_{i} > D\}|.$$

Similarly, the metrics for a category *C* is the sum of the metrics for all blocks in *C*:

$$S_C^{\text{AOM}}(D) = \sum_{b \in C} S_b^{\text{AOM}}(D), \quad U_C^{\text{AOM}}(D) = \sum_{b \in C} U_b^{\text{AOM}}(D), \quad W_C^{\text{AOM}}(D) = \sum_{b \in C} W_b^{\text{AOM}}(D).$$

The only difference between AdmitOnWrite and AdmitOnMiss is that AdmitOnWrite also takes into account write accesses. Therefore, for AdmitOnWrite, we let  $t_1$  be the write access time,  $t_2$ be the first read access time,  $t_3$  be the second read access time and so on. Then, we can similarly define  $S_C^{AOW}(D)$ ,  $U_C^{AOW}(D)$ , and  $W_C^{AOW}(D)$ .

For AdmitOnSecondMiss, a block is admitted at the second miss (read access). In addition, to prevent the cache from inserting a cold block, we require that, when inserting a block, its last read access time not be older than the oldest last access time of the blocks in the cache. Mathematically, a block is inserted at  $t_{i-1}$  (if not already in the cache) only if  $d_{i-1} = t_{i-1} - t_{i-2} \leq D$ . Therefore, the condition that a block is in the cache at  $t_{i-1}$ , either because it is already in the cache or it is

<sup>&</sup>lt;sup>3</sup>Bytes of occupied cache multiplied by seconds of residence time in cache. The same concept is also used in LHD [4].

inserted, is  $d_{i-1} \leq D$ , and so an access at  $t_i$  is a cache hit if and only if  $d_{i-1} \leq D$  and  $d_i \leq D$ :

$$B_b^{AOSM}(D, i) = \begin{cases} 1, & \text{Buffer Cache Hit at } t_i, \text{ using AOSM} \\ 0, & \text{Buffer Cache Miss at } t_i, \text{ using AOSM} \\ S_b^{AOSM}(D) = |\{i : \max(d_{i-1}, d_i) > D, B_b^{AOSM}(D, i) = 0\}|. \end{cases}$$

For  $U_b^{AOSM}(D)$ , the access at  $t_i$  contributes cache usage if either it is a cache hit,  $\max(d_i, d_{i-1}) \leq D$ , with residence time  $d_i$ , or a block insertion,  $d_i \leq D < d_{i-1}$ , with residence time D:

$$U_b^{\text{AOSM}}(D) = \text{Size}(b) \times \sum_i \left[ d_i \times \mathbf{1}_{\{\max(d_i, d_{i-1}) \le D\}} + D \times \mathbf{1}_{\{d_i \le D < d_{i-1}\}} \right]$$

 $W_h^{AOSM}(D)$  is the block size times the number of insertions:

$$W_b^{\text{AOSM}}(D) = \text{Size}(b) \times |\{i : d_i \le D < d_{i-1}\}|.$$

Of course,  $S_C^{AOSM}(D)$ ,  $U_C^{AOSM}(D)$ , and  $W_C^{AOSM}(D)$  can be defined similarly.

Because NeverAdmit does not insert any blocks at all,  $U_C^{NA}(D) = 0$ ,  $W_C^{NA}(D) = 0$ , and  $S_C^{NA}(D)$  is the number of buffer cache misses because of the accesses to the blocks in C:

$$\begin{split} B_b^{\mathrm{NA}}(D,i) &= \begin{cases} 1, & \text{Buffer Cache Hit at } t_i, \text{ using NA} \\ 0, & \text{Buffer Cache Miss at } t_i, \text{ using NA} \end{cases}, \\ S_b^{\mathrm{NA}}(D) &= |\{i: B_b^{\mathrm{NA}}(D,i) = 0\}|, \\ S_C^{\mathrm{NA}}(D) &= \sum_{b \in C} S_b^{\mathrm{NA}}(D). \end{split}$$

For the prefetch admissions PrefetchOnMiss, PrefetchOnSecondMiss, and PrefetchOnWrite, the difference is that the prefetch admissions view a prefetch region as a single, aligned block (except the write accesses for PrefetchOnWrite), and the metrics estimations are based on the inter-arrival times for prefetch regions, which can be calculated from the access times for the individual blocks in the prefetch region.

#### 9.3 Linear Program

We minimize the total cost by formulating a linear program. The cost function is the sum of the cost of disk reads and the cost of written bytes<sup>4</sup>:

$$V_C^p(D) = \text{Cost of } S_C^p(D) + \text{Cost of } W_C^p(D),$$

for  $p \in \{AOM, AOW, AOSM, POM, POW, POSM, NA\}$  and a given category *C*.

A category can receive *fractional* admission policies. For example, CacheSack may decide that it is optimal to apply AOM, AOW, AOSM, POM, POW, POSM, NA to 30%, 20%, 10%, 15%, 5%, 7%, 13% of blocks in *C*, respectively. Then, we can formulate a linear program that finds optimal policy fractions  $\alpha_C^p$ ,  $p \in \{AOM, AOW, AOSM, POM, POSM, NA\}$  to minimize the overall cost:

$$\min_{\alpha_{C}^{p}} \sum_{C} \sum_{p \in \mathbf{P}} \alpha_{C}^{p} V_{C}^{p}(D), \quad 0 \le \alpha_{C}^{p} \le 1, \quad \sum_{p \in \mathbf{P}} \alpha_{C}^{p} = 1, \quad \mathbf{P} = \{\text{AOM, AOW, AOSM, POM, POW, POSM, NA}\},$$
(1)

<sup>&</sup>lt;sup>4</sup>The cost of flash storage is fixed and hence is omitted in the objective function; other costs such as CPU, RAM, power, and network are a very small fraction of the TCO and are ignored.

ACM Transactions on Storage, Vol. 19, No. 2, Article 13. Publication date: March 2023.



Fig. 22. Example of Andrew's monotone chain convex hull algorithm applied to the admission policies. In this example, the solid lines (NA-AOSM, AOSM-POM, and POM-AOW) will be chosen. The dashed line (AOW-POW) is also a part of the lower convex hull, but we do not choose it, because its slope is positive.

subject to the capacity constraint that the cache usage should not exceed the given cache capacity  $U_{\text{total}}$ :

$$\sum_{C} \sum_{p \in \mathbf{P}} \alpha_{C}^{p} U_{C}^{p}(D) \leq U_{\text{total}}, \quad \mathbf{P} = \{\text{AOM, AOW, AOSM, POM, POW, POSM, NA}\}.$$

We note that if the LRU cache is perfectly modeled by the approach in Section 9.2, the resulting cache retention time of the LRU cache is exactly *D* after applying the optimal policy fractions.

# 9.4 Greedy Algorithm

Although the linear program in Equation (1) can be solved by a standard solver, we are able to solve it by a greedy algorithm with a simple transformation. It is especially beneficial for the production deployment because of the low-overhead and stability of the greedy algorithm, compared to a generic solver. We first note that the difference between the above linear program and a fractional knapsack problem [14] is that for each category, we need to decide *coupled* seven fractions (six degrees of freedom), instead of two fractions (one degree of freedom) in a fractional knapsack problem. Thus, the greedy algorithm in Reference [14] cannot be directly applied. However, we can use Andrew's lower convex hull algorithm [2] to decouple the dependency.

For a given category C, the lower convex hull formed by the points

$$\left\{ \left( U_{C}^{p}, V_{C}^{p} \right), p \in \left\{ \text{AOM, AOW, AOSM, POM, POW, POSM, NA} \right\} \right\}$$

is the lowest cost of *C* that can be generated among all convex combinations of the policies. For example, Figure 22 is the lower convex hull constructed by the given admission policies by using Andrew's algorithm. Let  $F_C(u)$  denote the lower convex hull formed above, as a mapping from cache usage *u* to the corresponding cost, for each category *C*. By dropping any line segments with non-negative slopes, all  $F_C$  are strictly decreasing, piecewise linear functions. Then, we can transform the linear program to a convex optimization problem:

$$\min_{u_C \ge 0} \sum_C F_C(u_C), \quad \sum_C u_C \le U_{\text{total}}.$$

We can then solve the above convex optimization problem by the steepest descent method (a greedy algorithm). We initialize  $u_C = 0$  for all *C* and iteratively decide each  $u_C$  as follows. We first

choose the line segment with the most negative slope among all line segments of  $F_C$  and change the value of the corresponding  $u_C$ . In the same fashion, we then choose the line segment with second most negative slope and change the value of the corresponding  $u_C$ , then the third most negative slope, and so on, until the sum of  $u_C$  reaches  $U_{\text{total}}$ .

Because we allow fractional policies, the category corresponding to the last chosen line segment generally has the optimal policy as a convex combination of two policies, and the optimal policy of any other category must be exactly one policy.

#### 9.5 Optimization Over Modeled Cache Retention Times

The linear program in Equation (1) is to find the optimal policy fractions for a *given* modeled cache retention time D, which cannot be known *a priori*. Thus, we need to solve the same optimization problem for *all* possible D:

$$\min_{D>0} \min_{\alpha_{C}^{p}} \sum_{C} \sum_{p \in \mathbf{P}} \alpha_{C}^{p} V_{C}^{p}(D), \quad 0 \leq \alpha_{C}^{p} \leq 1, \quad \sum_{p \in \mathbf{P}} \alpha_{C}^{p} = 1, \quad \mathbf{P} = \{\text{AOM, AOW, AOSM, POM, POW, POSM, NA}\},$$

subject to the same capacity constraint:

$$\sum_{C} \sum_{p \in \mathbf{P}} \alpha_{C}^{p} U_{C}^{p}(D) \le U_{\text{total}}, \quad \mathbf{P} = \{\text{AOM, AOW, AOSM, POM, POW, POSM, NA}\}$$

To do this, we can use a standard scalar-variable optimization approach like Brent's method [10] for  $0 < D \leq \mathcal{D}$ , where  $\mathcal{D}$  is a suitable upper bound. A brute-force approach may be even more practical for implementation: We simply solve the optimization problem for a set of reasonable retention times:  $0 < D_1 < D_2 < \cdots < D_m = \mathcal{D}$ .

#### 10 RELATED WORK

### **Production Flash Cache Algorithms**

Both Google [1] and Facebook [18] use FIFO-based evictions in their production flash caches to trade cache performance for reduced write amplification. RIPQ [36] is a non-FIFO, advanced flash cache algorithm that brings higher hit ratios while write amplification is well-controlled. Flashield [16] further improves RIPQ's write amplification by using DRAM as a buffer, and only writes flash-worthy objects into flash, predicted by a lightweight support vector machine classifier. CacheLib [7] resolved Flashield's issue that the TTLs of objects in the DRAM buffer are too short to be effective. CacheLib uses Bloom filters to count the number of accesses per object in the past 6 h (similar to TinyLFU [15]), to predict the number of accesses in the future, and uses FIFO for eviction. Kangaroo [26] further improves CacheLib's performance for tiny objects. DSS [28] uses predefined rules to classify I/O requests into different priorities, and applies heuristic admission and eviction policies to different priorities. DSS has been implemented in Intel's Cache Acceleration Software. Amazon's AQUA [3] analyzes workload patterns to place data into the appropriate tier.

CacheSack's high-level idea is similar to Flashield and CacheLib: keep the eviction simple to control write amplification, and use more sophisticated admission to improve cache performance and flash write endurance. For evictions, Flashield uses the CLOCK [12] approach and CacheSack uses Second Chance [30] to achieve LRU-style evictions. On the admission side, instead of using DRAM as a buffer, CacheSack has no in-memory buffer and expands the metadata table (ghost cache) for a more complete history; the median of the ghost cache TTL is 20 h, which is several times longer than the information used by Flashield and CacheLib. With a more complete history, CacheSack is able to build a more sophisticated model for admission. CacheSack also considers the two major costs of operating flash caches, disk reads and flash wearout, as a whole, and minimizes TCO.

CacheSack also utilizes the advantage that the categories are well-defined in the database systems served by Colossus Flash Cache. Classifying unstructured data is usually a difficult problem in machine learning. For Google's database systems, the classification is naturally available, and the categories often hint their cacheability.

#### **Admission Algorithms**

LARC [21] was previously used by Colossus Flash Cache as the default admission policy. LARC is designed for flash caches, and reduces write rate by inserting an object into the cache only when it is read a second time, based on the observation that most objects are read only once. Thus, inserting only the objects that are read a second time into the cache significantly reduces the write rate and the cache pollution. This strategy is particularly useful when a significant portion of the traffic is accessed only once, for example, Tencent's photo traffic [41] and AliCloud [24]. However, LARC loses all the second-access hits and becomes undesirable for long-tail accesses like Facebook's cache for social network photos. In the past, Colossus Flash Cache manually disabled LARC for workloads in which LARC underperformed. Selective admissions like TinyLFU [15] (non-window version) and HEC [42] that sacrifice the first few hits to determine the cacheability of data likely have the same issue.

TinyLFU [15] works by comparing the expected hit ratio of a newly accessed object against that of the object that would be evicted next from the cache, inserting the new object if its likely hit ratio is higher. Any eviction policy can be used to select the eviction victim (LRU is typical). TinyLFU predicts hit ratios for the objects using approximate counting (Bloom filters) of access frequency. TinyLFU also needs some extra structures to work properly: *Doorkeeper* is used to filter one-accessed blocks (the same use of LARC's ghost cache), and a DRAM buffer cache in front of TinyLFU (W-TinyLFU). All these structures require extra parameter tuning, which does not best fit the needs of Colossus Flash Cache as a general-purpose cache. mARC [32] uses ARC [27] as the eviction policy and dynamically determines whether to admit data on the first miss (naive ARC) or second miss (LARC). UBC [31] proposes a low-overhead mechanism to partition shared on-chip cache.

## **Eviction Algorithms**

There are also extensive studies on advanced eviction algorithms. Beckmann and Sanchez's method [5] evicts a block based on the block's economic value added. Instead of LRU or LFU that require specific data structures, Hyperbolic Caching [9] evicts a block based on a time-decay (hyperbolic) value function and uses a sampling technique to resolve the issue of the data structure requirement. Similarly, LHD [4] evicts the block of the lowest hit density, the number of hits per cache byte-second, and also applies a sampling technique to overcome the data structure issue. Hawkeye [22] assumes that the recent history can predict the near future and hence one can train a predictor learned from Belady's OPT [6] running on the recent traces. Waldspurger et al. [40] consider an ensemble of candidates, which can be a set of existing algorithms, or the same algorithm with different parameters, run scaled-down simulations on each candidate, and periodically adopt the most performant one.

#### **Machine Learning Algorithms**

With the recent rapid development of ML, there are also a few papers that adopt ML techniques to enhance cache performance. LFO [8] and LRB [35] use ML models to learn Belady's OPT [6], and apply the ML models to Content Delivery Networks caches. Parrot [25] also use ML to learn Belady's OPT from history, but uses modern deep learning architectures like Transformer [38] and BiDAF [33]. Reference [41] utilizes a concept similar to LARC [21] that the majority of traffic is

accessed just once, and uses ML models to predict whether data is worth inserting into the flash cache. The algorithm showed a large flash write reduction in Tencent's photo cache system as well as the improvement of hit ratios and latency. LeCaR [39] uses an ML approach to adaptively decide the better policy between LRU and LFU at eviction time. Zhou and Maas [43] model the inter-arrival times of a block as a log-normal distribution and learn the parameters from traces; then the evictions are executed in the manner of Belady's OPT: the block with lowest probability to get the next access in the near future will be evicted.

# 11 CONCLUSIONS

In this article, we introduce CacheSack, an admission policy optimization for Google's datacenter flash caches. CacheSack provides an efficient estimation for the performance metrics of an LRU-style cache under various configuration options. We use a knapsack approach to identify the optimal admission polices to minimize TCO. We share the experience of deploying CacheSack in Colossus Flash Cache, the general-purpose flash cache serving Colossus, which has since become the default admission policy. CacheSack requires less manual configuration than the previously used cache admission algorithm (LARC), significantly reduces disk reads and bytes written to flash, and improves TCO by 7.7% compared to LARC.

# ACKNOWLEDGMENTS

The authors thank Cory Casper, Martin Maas, and Richard McDougall for their assistance in various stages of the design and implementation of this algorithm. We also thank Dan Gibson, Larry Greenfield, Aaron Laursen, Milo Martin, Damodharan Rajalingam, and John Wilkes for their reviews and suggestions that improved this manuscript.

# REFERENCES

- [1] Christoph Albrecht, Arif Merchant, Murray Stokely, Muhammad Waliji, François Labelle, Nate Coehlo, Xudong Shi, and C. Eric Schrock. 2013. Janus: Optimal flash provisioning for cloud storage workloads. In *Proceedings of the USENIX Annual Technical Conference (USENIX ATC'13)*. USENIX Association, San Jose, CA, 91–102. Retrieved from https: //www.usenix.org/conference/atc13/technical-sessions/presentation/albrecht.
- [2] A.M. Andrew. 1979. Another efficient algorithm for convex hulls in two dimensions. Inform. Process. Lett. 9, 5 (1979), 216–219. https://doi.org/10.1016/0020-0190(79)90072-3
- [3] Jeff Barr. 2021. AQUA (Advanced Query Accelerator)—A Speed Boost for Your Amazon Redshift Queries. (April 2021). Retrieved from https://aws.amazon.com/blogs/aws/new-aqua-advanced-query-accelerator-for-amazon-redshift/.
- [4] Nathan Beckmann, Haoxian Chen, and Asaf Cidon. 2018. LHD: Improving cache hit rate by maximizing hit density. In Proceedings of the 15th USENIX Symposium on Networked Systems Design and Implementation (NSDI'18). USENIX Association, 389–403. Retrieved from https://www.usenix.org/conference/nsdi18/presentation/beckmann.
- [5] N. Beckmann and D. Sanchez. 2017. Maximizing cache performance under uncertainty. In Proceedings of the IEEE International Symposium on High-Performance Computer Architecture (HPCA'17). IEEE Computer Society, 109–120. https://doi.org/10.1109/HPCA.2017.43
- [6] L. A. Belady. 1966. A study of replacement algorithms for a virtual-storage computer. IBM Syst. J. 5, 2 (June 1966), 78–101. https://doi.org/10.1147/sj.52.0078
- [7] Benjamin Berg, Daniel S. Berger, Sara McAllister, Isaac Grosof, Sathya Gunasekar, Jimmy Lu, Michael Uhlar, Jim Carrig, Nathan Beckmann, Mor Harchol-Balter, and Gregory R. Ganger. 2020. The CacheLib caching engine: Design and experiences at scale. In *Proceedings of the 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI'20)*. USENIX Association, 753–768. Retrieved from https://www.usenix.org/conference/osdi20/presentation/berg.
- [8] Daniel S. Berger. 2018. Towards lightweight and robust machine learning for CDN caching. In Proceedings of the 17th ACM Workshop on Hot Topics in Networks (HotNets'18). ACM, 134–140. https://doi.org/10.1145/3286062.3286082
- [9] Aaron Blankstein, Siddhartha Sen, and Michael J. Freedman. 2017. Hyperbolic caching: Flexible caching for web applications. In *Proceedings of the USENIX Annual Technical Conference (USENIX ATC'17)*. USENIX Association, 499–511. Retrieved from https://www.usenix.org/conference/atc17/technical-sessions/presentation/blankstein.
- [10] R. P. Brent. 1973. Algorithms for Minimization Without Derivatives. Prentice-Hall.

#### CacheSack: Theory and Experience of Google's Admission Optimization

- [11] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber. 2008. Bigtable: A distributed storage system for structured data. ACM Trans. Comput. Syst. 26, 2, Article 4 (June 2008), 26 pages. https://doi.org/10.1145/1365815.1365816
- [12] F. J. Corbató. 1968. A Paging Experiment with the Multics System. Massachusetts Institute of Technology. Retrieved from https://books.google.com/books?id=5wDQNwAACAAJ.
- [13] James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, J. J. Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, Wilson Hsieh, Sebastian Kanthak, Eugene Kogan, Hongyi Li, Alexander Lloyd, Sergey Melnik, David Mwaura, David Nagle, Sean Quinlan, Rajesh Rao, Lindsay Rolig, Yasushi Saito, Michal Szymaniak, Christopher Taylor, Ruth Wang, and Dale Woodford. 2012. Spanner: Google's globally distributed database. In Proceedings of the 10th USENIX Symposium on Operating Systems Design and Implementation (OSDI'12). USENIX Association, 261–264. Retrieved from https://www.usenix.org/conference/osdi12/technicalsessions/presentation/corbett.
- [14] George B. Dantzig. 1957. Discrete-variable extremum problems. Oper. Res. 5, 2 (1957), 266–288. https://doi.org/10. 1287/opre.5.2.266. arXiv:https://doi.org/10.1287/opre.5.2.266
- [15] Gil Einziger, Roy Friedman, and Ben Manes. 2017. TinyLFU: A highly efficient cache admission policy. ACM Trans. Storage 13, 4, Article 35 (Nov. 2017), 31 pages. https://doi.org/10.1145/3149371
- [16] Assaf Eisenman, Asaf Cidon, Evgenya Pergament, Or Haimovich, Ryan Stutsman, Mohammad Alizadeh, and Sachin Katti. 2019. Flashield: A hybrid key-value cache that controls flash write amplification. In *Proceedings of the 16th USENIX Symposium on Networked Systems Design and Implementation (NSDI'19)*. USENIX Association, 65–78. Retrieved from https://www.usenix.org/conference/nsdi19/presentation/eisenman.
- [17] Ronald Fagin. 1977. Asymptotic miss ratios over independent references. J. Comput. Syst. Sci. 14, 2 (1977), 222–250. https://doi.org/10.1016/S0022-0000(77)80014-7
- [18] Alex Gartrell. 2013. McDipper: A key-value cache for Flash storage. (March 2013). Retrieved from https://engineering. fb.com/2013/03/05/web/mcdipper-a-key-value-cache-for-flash-storage/.
- [19] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. 2003. The Google file system. In Proceedings of the 19th ACM Symposium on Operating Systems Principles. ACM, 20–43.
- [20] Dean Hildebrand and Denis Serenyi. 2021. Colossus under the hood: A peek into Google's scalable storage system. Retrieved from https://cloud.google.com/blog/products/storage-data-transfer/a-peek-behind-colossus-googles-filesystem.
- [21] Sai Huang, Qingsong Wei, Dan Feng, Jianxi Chen, and Cheng Chen. 2016. Improving flash-based disk cache with lazy adaptive replacement. ACM Trans. Stor. 12, 2, Article 8 (Feb. 2016), 24 pages. https://doi.org/10.1145/2737832
- [22] Akanksha Jain and Calvin Lin. 2016. Back to the future: Leveraging Belady's algorithm for improved cache replacement. In Proceedings of the 43rd International Symposium on Computer Architecture (ISCA'16). IEEE Press, 78–89. https://doi.org/10.1109/ISCA.2016.17
- [23] Bo Jiang, Philippe Nain, and Don Towsley. 2018. On the convergence of the TTL approximation for an LRU cache under independent stationary request processes. ACM Trans. Model. Perform. Eval. Comput. Syst. 3, 4, Article 20 (Sept. 2018), 31 pages. https://doi.org/10.1145/3239164
- [24] Jinhong Li, Qiuping Wang, Patrick P. C. Lee, and Chao Shi. 2020. An in-depth analysis of cloud block storage workloads in large-scale production. In *Proceedings of the IEEE International Symposium on Workload Characterization* (*IISWC*'20). IEEE Computer Society, 37–47. https://doi.org/10.1109/IISWC50251.2020.00013
- [25] Evan Zheran Liu, Milad Hashemi, Kevin Swersky, Parthasarathy Ranganathan, and Junwhan Ahn. 2020. An imitation learning approach for cache replacement. In *Proceedings of the 37th International Conference on Machine Learning* (ICML'20) (Proceedings of Machine Learning Research), Vol. 119. PMLR, 6237–6247. Retrieved from http://proceedings. mlr.press/v119/liu20f.html.
- [26] Sara McAllister, Benjamin Berg, Julian Tutuncu-Macias, Juncheng Yang, Sathya Gunasekar, Jimmy Lu, Daniel S. Berger, Nathan Beckmann, and Gregory R. Ganger. 2021. Kangaroo: Caching billions of tiny objects on flash. In Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles (SOSP'21). ACM, 243–262. https: //doi.org/10.1145/3477132.3483568
- [27] Nimrod Megiddo and Dharmendra S. Modha. 2003. ARC: A self-tuning, low overhead replacement cache. In Proceedings of the 2nd USENIX Conference on File and Storage Technologies (FAST 03). USENIX Association. Retrieved from https://www.usenix.org/conference/fast-03/arc-self-tuning-low-overhead-replacement-cache.
- [28] Michael Mesnier, Feng Chen, Tian Luo, and Jason B. Akers. 2011. Differentiated storage services. In Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP'11). ACM, 57–70. https://doi.org/10.1145/2043556. 2043563
- [29] Changwoo Min, Kangnyeon Kim, Hyunjin Cho, Sang-Won Lee, and Young Ik Eom. 2012. SFS: Random write considered harmful in solid state drives. In *Proceedings of the 10th USENIX Conference on File and Storage Technologies* (FAST'12). USENIX Association, 12.

- [30] Pancham Pancham, Deepak Chaudhary, and Ruchin Gupta. 2014. Comparison of cache page replacement techniques to enhance cache memory performance. *Int. J. Comput. Appl.* 98 (July 2014), 27–33. https://doi.org/10.5120/17293-7771
- [31] Moinuddin K. Qureshi and Yale N. Patt. 2006. Utility-based cache partitioning: A low-overhead, high-performance, runtime mechanism to partition shared caches. In Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO'06). 423–432. https://doi.org/10.1109/MICRO.2006.49
- [32] Ricardo Santana, Steven Lyons, Ricardo Koller, Raju Rangaswami, and Jason Liu. 2015. To ARC or not to ARC. In Proceedings of the 7th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage'15). USENIX Association. Retrieved from https://www.usenix.org/conference/hotstorage15/workshop-program/presentation/santana.
- [33] Min Joon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2017. Bidirectional attention flow for machine comprehension. In *Proceedings of the 5th International Conference on Learning Representations (ICLR'17)*. OpenReview.net. Retrieved from https://openreview.net/forum?id=HJ0UKP9ge.
- [34] Jeff Shute, Radek Vingralek, Bart Samwel, Ben Handy, Chad Whipkey, Eric Rollins, Mircea Oancea, Kyle Littlefield, David Menestrina, Stephan Ellner, John Cieslewicz, Ian Rae, Traian Stancescu, and Himani Apte. 2013. F1: A distributed SQL database that scales. *Proc. VLDB Endow.* 6, 11 (2013).
- [35] Zhenyu Song, Daniel S. Berger, Kai Li, and Wyatt Lloyd. 2020. Learning relaxed belady for content distribution network caching. In *Proceedings of the 17th USENIX Symposium on Networked Systems Design and Implementation* (NSDI'20). USENIX Association, 529–544. Retrieved from https://www.usenix.org/conference/nsdi20/presentation/ song.
- [36] Linpeng Tang, Qi Huang, Wyatt Lloyd, Sanjeev Kumar, and Kai Li. 2015. RIPQ: Advanced photo caching on flash for Facebook. In Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'15). USENIX Association, 373–386. Retrieved from https://www.usenix.org/conference/fast15/technical-sessions/presentation/tang.
- [37] Rajesh Thallam. 2020. BigQuery explained: An overview of BigQuery's architecture. Retrieved from https://cloud. google.com/blog/products/data-analytics/new-blog-series-bigquery-explained-overview.
- [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Retrieved from https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- [39] Giuseppe Vietri, Liana V. Rodriguez, Wendy A. Martinez, Steven Lyons, Jason Liu, Raju Rangaswami, Ming Zhao, and Giri Narasimhan. 2018. Driving cache replacement with ML-based LeCaR. In Proceedings of the 10th USENIX Conference on Hot Topics in Storage and File Systems (HotStorage'18). USENIX Association, 3.
- [40] Carl Waldspurger, Trausti Saemundsson, Irfan Ahmad, and Nohhyun Park. 2017. Cache modeling and optimization using miniature simulations. In *Proceedings of the USENIX Annual Technical Conference (USENIX ATC'17)*. USENIX Association, 487–498. Retrieved from https://www.usenix.org/conference/atc17/technical-sessions/presentation/ waldspurger.
- [41] Hua Wang, Xinbo Yi, Ping Huang, Bin Cheng, and Ke Zhou. 2018. Efficient SSD caching by avoiding unnecessary writes using machine learning. In Proceedings of the 47th International Conference on Parallel Processing (ICPP'18). ACM, Article 82, 10 pages. https://doi.org/10.1145/3225058.3225126
- [42] Jingpei Yang, Ned Plasson, Greg Gillis, Nisha Talagala, Swaminathan Sundararaman, and Robert Wood. 2013. HEC: Improving endurance of high performance flash-based cache devices. In *Proceedings of the 6th International Systems and Storage Conference (SYSTOR'13)*. ACM, Article 10, 11 pages. https://doi.org/10.1145/2485732.2485743
- [43] Giulio Zhou and Martin Maas. 2021. Learning on distributed traces for data center storage systems. In *Proceedings of Machine Learning and Systems*, A. Smola, A. Dimakis, and I. Stoica (Eds.), Vol. 3. 350–364. Retrieved from https: //proceedings.mlsys.org/paper/2021/file/82161242827b703e6acf9c726942a1e4-Paper.pdf.
- [44] Huapeng Zhou, Linpeng Tang, Qi Huang, and Wyatt Lloyd. 2016. The Evolution of Advanced Caching in the Facebook CDN. Retrieved from https://research.fb.com/blog/2016/04/the-evolution-of-advanced-caching-in-thefacebook-cdn/.

Received 20 December 2022; accepted 9 January 2023

#### 13:24