

Lightweight Robust Size Aware Cache Management

Gil Einziger
Computer Science Department
Ben Gurion University
gilga1983@gmail.com

Ohad Eytan
Computer Science Department
Technion
ohadey@cs.technion.ac.il

Roy Friedman
Computer Science Department
Technion
roy@cs.technion.ac.il

Benjamin Manes
Independent
ben.manes@gmail.com

Abstract

Modern key-value stores, object stores, Internet proxy caches, as well as Content Delivery Networks (CDN) often manage objects of diverse sizes, e.g., blobs, video files of different lengths, images with varying resolution, and small documents. In such workloads, size-aware cache policies outperform size-oblivious algorithms. Unfortunately, existing size-aware algorithms tend to be overly complicated and computationally expensive.

Our work follows a more approachable pattern; we extend the prevalent (size-oblivious) TinyLFU cache admission policy to handle variable sized items. Implementing our approach inside two popular caching libraries only requires minor changes. We show that our algorithms yield competitive or better hit-ratios and byte hit-ratios compared to the state of the art size-aware algorithms such as AdaptSize, LHD, LRB, and GDSF. Further, a runtime comparison indicates that our implementation is faster by up to x3 compared to the best alternative, i.e., it imposes much lower CPU overhead.

1 Introduction

Caches improve the performance of diverse computing and networked systems by storing portions of the data closer to the applications and users. Caching is widely used in many databases and storage systems in which DRAM memory serves as a cache for slower secondary storage [25], Web caching [34], and even in hardware fast SRAM memory serves as a cache to slower DRAM memory. As another example, *Content Delivery Networks* (CDN) reduce the data access latency by storing popular content closer (in terms of Internet topology) to its consumers [22]. CDNs deploy multiple distributed content caching servers, also known as *Points of Presence* (PoP), in multiple ASs. PoPs serve popular content faster than fetching the data from the original server and accomplish three goals: reducing Internet traffic [21], reducing the load on the content providers' servers, and reducing users' perceived response times, thereby improving the user experience.

Given that the cache is smaller than the entire storage whose content it stores, a *cache management policy* decides which items to store in the cache at any given point. Due to the diversity of workloads and systems, a plethora of such policies have been devised, e.g., [27, 33, 10, 43, 19, 5, 6, 41, 24, 14, 15, 45].

The vast majority of research on cache management policies ignore the size of cached items. The logic behind this simplification is that many databases and storage systems use block devices as secondary storage, so caches only need to work at block boundary. The sizes of all blocks or pages within a given system are the same or very similar.

However, in modern databases and data-stores, the reality is more complicated. Specifically, there are many important scenarios, such as imaging and video objects, CDN's caches, cell content in column-family data models, blobs, and application caches, where the size of cache-able items vary, emphasizing the need for

size-aware cache policies. Further, much of today’s Internet content volume is consumed by video streaming, which is encoded and transmitted using variable bit rates. Therefore, the video is partitioned into chunks (each corresponding to a short video duration) whose size varies considerably [3]. In fact, even the same sub-stream may be encoded multiple times at different bit rates, each of which is transmitted to different clients based on their temporal line quality. In such workloads, size-aware cache policies often outperform size oblivious cache policies. Indeed, a modern analysis has recently shown that sophisticated variable sized objects caching indeed has room for significant system performance improvements [9].

Alas, size-aware cache management policies are often more complex than unit sized policies. For example, size-aware policies may divide the cache into *slabs*, where each slab is devoted to objects of similar size and is managed independently of other slabs [39, 6]. Such an approach leads to sub-optimal cache utilization as objects of a certain slab may be more popular than objects of other slabs. Further, the popularity of objects varies, so there is no single best static slabbing partitioning. Some works mitigate these complications by periodically re-balancing the slab allocation [6].

When objects’ sizes vary, we distinguish between *hit-ratio* and *byte hit-ratio*. Hit-ratio is the ratio between the number of accesses that are a hit in the cache and the total number of accesses. Hit-ratio is important to improve the systems’ perceived access latency, and thereby its user experience. In contrast, byte hit-ratio is the accumulated number of bytes served by hits in the cache vs. the total number of bytes returned by all accesses. This indicates the cache’s effectiveness in reducing network bandwidth consumption.

The seminal GDSF policy [13] is very effective in terms of both hit ratio and byte hit-ratio. But, GDSF suffers from a poor reputation of being computationally expensive and it operates in logarithmic complexity. Three other prominent size-aware cache management policies include AdaptSize [10], LHD [6] and the recent LRB [41]. AdaptSize involves parameter tuning through non-trivial stochastic models. LHD and LRB both apply a heavy machine learning process to learn from previous accesses to predict future ones. Here we are looking for a lightweight, yet effective, size aware caching policy.

Our Contributions

This paper explores the benefits of simple size-aware cache policies that generalize the W-TinyLFU policy [20] to handle varying sized objects¹. In a gist, W-TinyLFU decides which objects to hold in the cache for a long duration according to their access frequency compared to the cache victim. Intuitively, variable-sized items imply that we may need to compare a newly arriving item against multiple cache victims to make enough room for it in the cache.

Our work evaluates three such possible extensions. The first two, nicknamed *Implicit Victims* (IV) and *Queue of Victims* (QV), are implemented in production extensions of W-TinyLFU to variable-sized objects (our Caffeine [32] and Ristretto [18]), but were never described outside their code, nor analyzed and compared in any manner. The third is our newly proposed approach, nicknamed *Aggregated Victims* (AV). We demonstrate that our proposed method AV outperforms the two existing solutions IV and QV in hit-ratio. When it comes to byte hit-ratio, QV has a slight advantage in the CDN traces. Next, we show that our new policy AV achieves competitive hit-ratio and byte hit-ratio compared to the state-of-the-art policies GDSF [13], AdaptSize [10], LHD [6] and LRB [41] over a wide variety of real-world traces. While none of the evaluated policies is best in all the cases, our policy consistently attains high hit ratios and faster constant time operation and is also very simple to implement.

Interestingly, our measurements revealed the following inherent limitation of cache management approaches, such as AdaptSize [10], that admit objects with a probability that is inversely proportional to the object’s size. The problem is that such policies fail to utilize the entire cache area when operating over large caches, as found in modern CDN PoPs [22]. We explain this problem in the evaluation section.

Paper road-map: The rest of this paper is structured as follows: We survey related work in Section 2. We present the (size oblivious) W-TinyLFU approach in Section 3, and the extensions to variable size objects

¹W-TinyLFU is used, directly or through caching libraries, in Cassandra, Accumulo, HBase, Apache Solr, Infinispan, OpenWhisk, Corfu, Finagle, Spring, Akka, Neo4j, DGraph, Druid, and many others as listed in [32].

in Section 4. Section 5 presents the performance evaluation of the hit-ratio, byte hit-ratio, and throughput for the various approaches. Finally, we conclude with a discussion in Section 6.

2 Related Work

As indicated above, the majority of cache management policies assume that all objects are of similar or equal size. Such policies include LFU and its variants [28, 29, 4] as well as LRU [23] and its performance-driven approximations [44, 24]. More modern approaches improve on LRU by adding frequency aspects to it. These include, e.g., LRU-K [37], 2Q [27], ARC [33], CAR [5], LIRS [25], DLIRS [31], FRD [38] and Hyperbolic [12]. Mini-Sim [45] efficiently simulates multiple policies and configurations in order to choose the best one. W-TinyLFU proposes having two cache segments with an approximated LFU admission filter for transferring items between these segments [19, 20].

The seminal Greedy-Dual work was one of the first to incorporate size into the replacement decision [46]. The Greedy-Dual was later extended to accommodate for recency, frequency, cost, and size in the GDSF scheme [13]. GDSF gives a dynamic score to each item by combining the above four criteria into a unifying formula. To make room for a newly arriving item x , GDSF collects cached items whose dynamic score is minimal, until finding enough items whose combined size is larger than the size of x and then evicts all these found items. To date, GDSF remains a tough performance benchmark to beat as far as hit ratio and byte hit ratio are concerned. Yet, its computational complexity is largely considered too high for most production caching systems.

The popular Redis system [39] divides the total cache area into multiple slabs, each dedicated to similarly-sized objects. Every slab is maintained using either sampled LRU or sampled LFU for eviction coupled with blind admission.

Three promising recent size aware caching schemes are AdaptSize [10], LHD [6] and LRB [41]. Adapt-Size implements a stochastic sized-based admission policy where large items have a smaller probability of being admitted compared with small items. Specifically, AdaptSize utilizes Markov models to optimize the admission policy.

In contrast, LHD [6] introduced the concept of *hit density* to capture the combination of frequency and size. It also considers recency through machine learning based prediction of when each item would be reaccessed. These parameters are placed in a unifying formula, and the cache victims are those whose score is minimal. Yet, to keep the computational cost practical, LHD randomly samples 64 uniformly selected items, computes the scores only for those items, and evicts the minimal among them. LHD also employs slabbing to avoid having to evict too many items at once, as well as slabs re-balancing to improve the overall cache utilization.

LRB [41] is a recent sophisticated cache replacement policy whose design goal is to maximize byte hit ratio in order to minimize bandwidth consumption in CDNs. LRB employs machine learning in order to approximate the behavior of a *relaxed Belady* algorithm. The classical offline Belady algorithm [8] evicts from the cache the object whose future access is the furthest, which is optimal. Yet, as the authors of [41] have discovered, learning to approximate Belady is very hard, but is also usually not needed. Instead, they proposed a variant called relaxed Belady, in which the algorithm evicts some object whose next access is beyond what they referred to as the *Belady boundary*. The LRB policy utilizes machine learning in order to mimic the relaxed Belady algorithm, which is feasible and indeed provides very good byte hit ratio results.

3 An Overview of W-TinyLFU

In this section, we give a brief overview of the (size oblivious) W-TinyLFU scheme. For more details, see [19, 20]. As shown in Figure 1, W-TinyLFU relies on three modules: the *Main cache*, the *Window cache* and an approximated LFU based *admission filter* called *TinyLFU*. On a cache miss, the accessed item is inserted into the Window cache. When this results in another item being evicted from Window cache, aka *Window victim*, the Window victim is introduced to the TinyLFU filter. The filter provides approximate

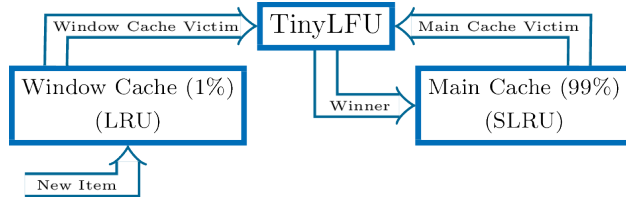


Figure 1: W-TinyLFU scheme: Items are first always admitted to the Window cache. The Window cache victim is offered to the Main cache, which employs TinyLFU as its admission filter. We preserve here the terminology of [20].

frequency estimations of the Window victim, and the would-be victim of the Main cache, aka *Main victim*. If the Window victim’s estimated frequency is higher than that of the Main victim, then the Window victim is admitted to the Main cache at the expense of the Main victim. Else, the Main victim remains in the Main cache, and the Window victim gets evicted.

In principle, this is a generic architecture where one can use any cache management scheme for the Main and Window caches, and the admission filter of the Main cache could employ any utility comparison function. Yet, the authors of [20] recommended deploying LRU as the eviction policy of the Window cache, SLRU as the eviction scheme of the Main cache, and TinyLFU as the admission filter. As for the relative sizes of the two cache regions, in [20] it was recommended that the Window cache be allocated 1% of the overall cache size, whereas a followup work [19] showed how to adapt these region sizes dynamically.

The TinyLFU admission filter is implemented through a sketch such as a *minimal increment counting Bloom filter* [16], or a *count min sketch* [17]. All sketch’ counters are halved for aging purposes every S accesses, where S is a parameter. Given a cache of size C items, these counters are also capped by S/C and therefore require $O(\log S/C)$ bits. Typically, S is at least one order of magnitude larger than C . The sketch counters corresponding to x are updated for every occurrence of x , even if it is not in the cache. Such an approach eliminates the need for maintaining ghost entries, as is common in other schemes [33, 26, 38].

4 Size-aware Admission Policy

This section discusses how to extend the W-TinyLFU policy to variable-sized objects. There are three main differences between the similar size case and variable size scenarios: First, a new item might be larger than the Window cache. In this case, the newly arriving item skips the Window cache and is immediately submitted for the TinyLFU filter to be determined if it can enter the Main cache.

Second, we may need to consider multiple potential Window victims whose aggregated size is large enough to make room for the newly arriving item. Figure 2 exemplifies this point. In this example, B is the new item, while W and V are the Window cache victims.

Third, in W-TinyLFU, the Window cache victim is tested against the Main cache victim. Here, we may have multiple Window cache victims, each tested against (possibly multiple) Main cache victims. There are several options on how to pick these potential Main cache victims and how to compare them to the candidate item(s). Specifically, as mentioned before, in this work, we consider three options. These include Implicit Victims (IV), Queue of Victims (QV), and Aggregated Victims (AV). The Caffeine Java-based project [32] employs IV, the Go-based Ristretto project [18] implements QV, and AV is our newly designed scheme. While IV and QV were adapted into popular libraries, we are the first to evaluate them systematically.

Figure 3 outlines the size-aware W-TinyLFU policies. The differences with the basic W-TinyLFU policy are marked in yellow. In our measurements, we use a Window cache whose size is 1% of the total cache size (and 99% is allocated to the Main cache), as suggested in [20]. The eviction policy for the Window cache is LRU, while for the Main cache, we experiment with few alternatives, as described below. Algorithm 1 provides pseudocode for miss handling by the size-aware adaptation of W-TinyLFU.

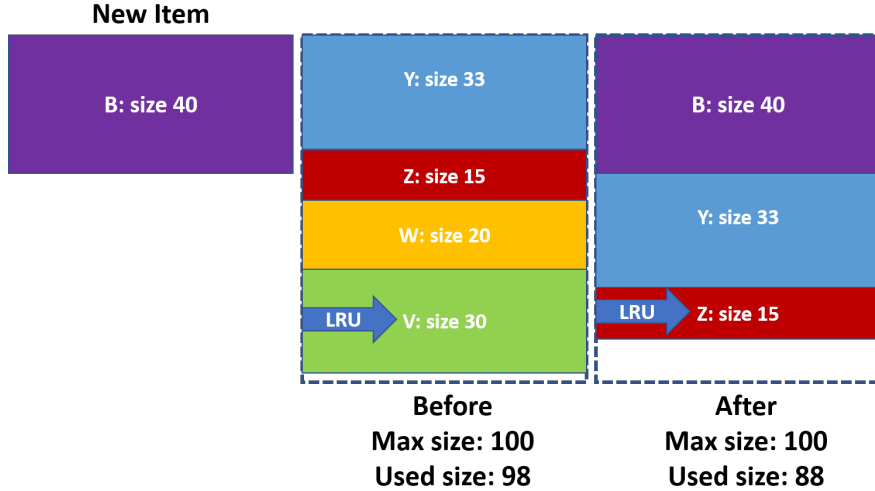


Figure 2: An illustration of the admission of a new item (B) to the Window cache, in this example, evicting the LRU item (V), does not make enough room for B, and we need to evict also the 2nd LRU item (W).

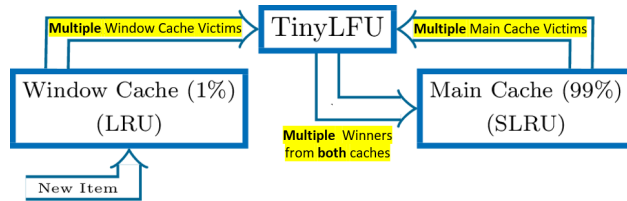


Figure 3: We are extending W-TinyLFU to size-awareness. The main change to the (size oblivious) W-TinyLFU policy is that adding the new item to the Window cache may result in multiple victims (in yellow). We use TinyLFU to decide for each one of these Window cache victims. Similarly, admitting a Window cache victim to the Main cache may result in multiple Main cache victims. Thus, in the end, the set of victims may come from both caches.

4.1 Implicit Victims (IV)

The IV scheme deals with variable sized objects by following the simplistic approach in Algorithm 2: IV compares the candidate item’s frequency vs. the frequency of the would-be victim, as reported by the approximate filter (Line 3). If the candidate is more frequent, IV blindly evicts enough victims, one after the other, until there is enough free space in the cache for the candidate to be admitted (Lines 4–6). The method is exemplified in Figure 4.

Notice that IV only cares about the first victim’s frequency and ignores the frequency of subsequent victims that may be evicted to admit the candidate. In case the victim is more frequent, IV promotes it (Line 9) and rejects the candidate (Line 10).

This approach is straightforward and computationally efficient to implement. As mentioned before, we developed IV and implemented it in Caffeine [32], but it was never explained before beyond the source code of Caffeine nor studied systematically.

Algorithm 1 W-TinyLFU Becomes Size Aware

```
1: procedure ONMISS(newItem)
2:   if newItem size > cache size then
3:     reject newItem
4:     return
5:   candidates  $\leftarrow$  emptyList()
6:   if newItem size > maximal window size then           ▷ The new item is too large for the Window cache
7:     candidates.append(newItem)                           ▷ Skip Window cache  $\rightarrow$  a candidate for Main
8:   else
9:     insert newItem to window
10:  while window size > maximal window size do           ▷ Evict items from Window to make room
11:    candidate  $\leftarrow$  evict from window                 ▷ Each eviction from Window becomes
12:    candidates.append(candidate)                         ▷ a candidate for Main
13:  for candidate in candidates do
14:    EvictOrAdmit(candidate)                               ▷ This function is where IV, QV and AV differ
```

Algorithm 2 Implicit Victims (IV)

```
1: procedure EVICTORADMIT(candidate)
2:   victim  $\leftarrow$  getVictim()
3:   if frequency(candidate)  $\geq$  frequency(victim) then
4:     while space is needed do
5:       victim  $\leftarrow$  getVictim()
6:       evict victim
7:       admit candidate
8:   else
9:     promote victim
10:  reject candidate
```

4.2 Queue of Victims (QV)

QV takes a somewhat different approach to handle variable-sized objects, as outlined in Algorithm 3. Specifically, QV repeatedly seeks for victims. For each victim, if the frequency of the candidate is higher than the victim’s frequency, the victim is evicted (Lines 4–5). Otherwise, the victim is promoted in the cache, as if it was accessed, and the iterative process stops (Lines 7–8). At the end of the iterative process, if there is enough vacant space in the cache for the candidate, then the candidate is admitted (Line 10).

Figure 5 illustrates this process. In this example, we evict an item (J) from the Main cache without admitting an item. The lack of admission is because the evicted item was not large enough to make room for the candidate, while the next victim was more frequent than the candidate.

The motivation behind QV, compared to IV, is two-fold: (i) avoid evicting items that are more frequent than the candidate, which could happen with IV, (ii) yet, if a victim is less frequent than the candidate, evict the victim even if eventually the candidate is not admitted. The rationale behind (ii) is that if a victim is less frequent than a candidate, it is polluting the cache, and so it should be evicted in any case.

4.3 Aggregated Victims (AV)

We now present our Aggregated Victims (AV) approach. AV compares a cache candidate’s frequency to that of the entire group of potential victims. We admit the candidate if its frequency is higher than that of *all* the potential cache victims. Intuitively, since admitting the candidate causes all cache victims’ eviction, the admission policy asserts that such a transaction is beneficial to the cache as a whole.

Algorithm 4 details AV’s behavior: Specifically, AV gathers potential victims (according to the eviction policy) until their total size is sufficient to admit the candidate (Lines 3–5). AV then uses TinyLFU to compare the candidates’ frequency to the total frequency of all victims (Line 8). If the candidate’s frequency

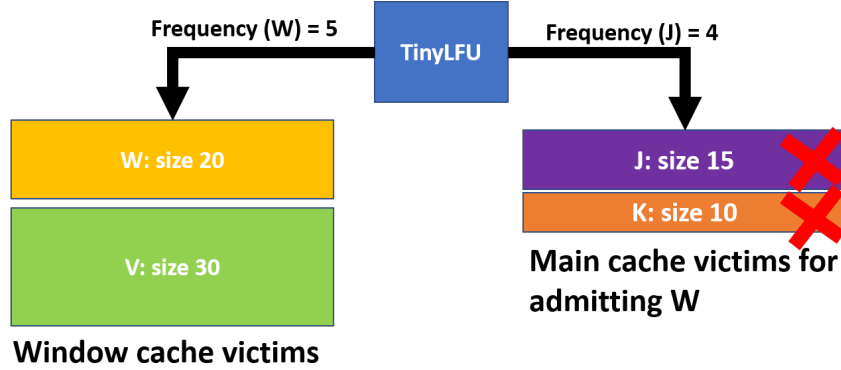


Figure 4: An example of the Implicit Victims (IV) approach: We compare the frequency of the Window cache victim (W) to that of the *first* Main cache victim (J). Since W is more frequent, we admit W and evict J and K .

Algorithm 3 Queue of Victims (QV)

```

1: procedure EVICTORADMIT(candidate)
2:   while space is needed do
3:     victim  $\leftarrow$  getVictim()
4:     if frequency(candidate)  $\geq$  frequency(victim) then
5:       evict victim
6:     else
7:       promote victim
8:       break
9:   if there is enough space then
10:    admit candidate
11:  else
12:    reject candidate

```

outweighs the total frequency of all victims, then it is admitted at the expense of the victims (Lines 9–11). Otherwise, the candidate is discarded, and no entries are evicted from the cache (Lines 13–15). As before, if the candidate is rejected, we promote the potential victims (Line 14). Figure 6 exemplifies this process. In the figure, W fails to enter the Main cache and is dropped from the cache entirely.

In most simple eviction policies (e.g., LRU), we also promote the victims as if they were accessed once after rejecting the candidate, to force the eviction policy to select a different set of victims for the next candidate. Yet, some eviction policies may not require this step.

4.3.1 The Early Pruning Optimization

The description of AV so far suffers from inefficiency. If the victims’ cumulative frequencies surpass the candidate’s frequency during its iterative process, we can abort the calculation since the candidate will not be admitted. Hence, we can improve run-time by gradually gathering the victims and summing their frequency estimations until the total victims’ frequency outweighs the candidate’s frequency, or their total size is sufficient to admit the candidate. This optimization is listed in (Lines 6–8) of Algorithm 4; we name it *early pruning*.

Early pruning improves the run time of AV since we need to compare rejected candidates (which are often the vast majority [2]) with fewer victims. Figure 7 demonstrates the effectiveness of this optimization. Notice that early pruning reduces victim comparisons by a factor of $x4-x16$.

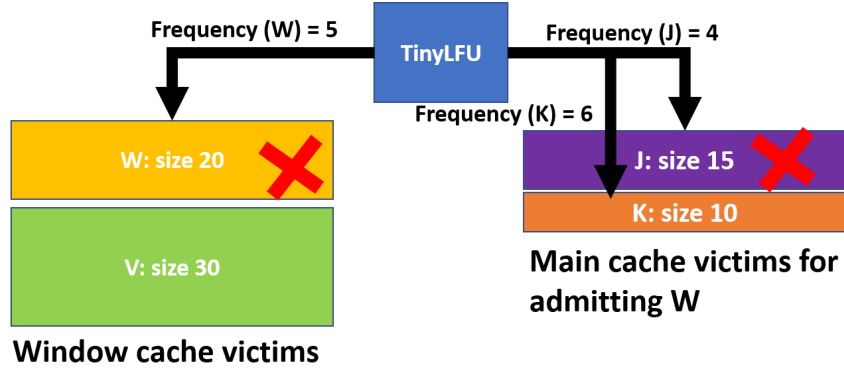


Figure 5: The Queue of Victims (QV) approach: Here, we compare the frequency of the Window cache victim (W) to that of the Main cache victim (J). Since W is more frequent, evict J and move to the next victim (K). In this case, K is more frequent than W, so we do not admit W to the Main cache.

Algorithm 4 Aggregated Victims (AV)

```

1: procedure EVICTORADMIT(candidate)
2:   victims  $\leftarrow$  emptyList()
3:   while victims size < space needed do
4:     victim  $\leftarrow$  getVictim()
5:     victims.append(victim)
6:     if frequency(candidate) < frequency(victims) then
7:       break
8:     if frequency(candidate)  $\geq$  frequency(victims) then
9:       for victim in victims do
10:        evict victim
11:        admit candidate
12:     else
13:       for victim in victims do
14:        promote victim
15:       reject candidate

```

\triangleright early pruning

Let us note that potential victims that were eventually not evicted are promoted in all these algorithms, e.g., Line 14 of Algorithm 4. Hence, early pruning may result in fewer objects being promoted than without this optimization. Yet, we empirically found that its impact on the hit ratio is negligible.

5 Evaluation

In this section, we present our conducted trace-based simulations results. First, we compare the three discussed approaches of the size-aware admission filter (IV, QV, and AV), combining six different Main cache eviction policies for W-TinyLFU: **SLRU** as in Caffeine [32], four versions of sampled policies and random eviction (**Random**). Each of the sampling policies samples five entries and selects the victim according to one of the following four rules respectively: lowest frequency (**Sampled Frequency**), largest size (**Sampled Size**), lowest frequency/size (**Sampled Frequency/Size**), and the closest match where the evicted item is the one yielding highest memory utilization (**Sampled Needed Size**), e.g., when admitting a 100k item into a cache with 50k free space, the victim is the one whose size is closest to 50k). Our sampled configurations mimic Ristretto’s SampledLFU configuration [18] that randomly samples five potential victims and evicts the one with the lowest frequency. Overall, we suggest eighteen combinations of admission and eviction

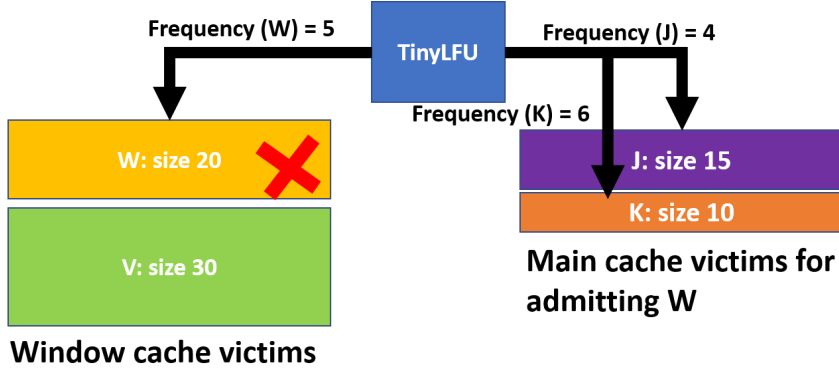


Figure 6: Illustrating the Aggregate Victims (AV) approach: In this example, we compare the frequency of the Window cache victim (W) to that of the Main cache victims (J, K). The combined frequency of the Main cache victims ($6+4=10$) is larger than the frequency of W (5). Therefore, we do not admit W to the Main cache.

policies in total.

Next, we picked our top versions, AV with SLRU for hit-ratio and QV with SLRU for byte-hit-ratio, and compared it to four other leading size-aware policies: GDSF [13], AdaptSize [10], LHD [6] and LRB [41]. Following previous works, our primary measurement was hit-ratio, byte-hit-ratio, and run-time performance. As mentioned before, hit-ratio improves latencies, resulting in better user experience and byte-hit-ratio indicates how much bandwidth is saved, meaning better network utilization, while run-time points to the CPU overhead for maintaining the cache itself.

To perform the evaluation, we used our own Java implementation to measure all but AdaptSize, LHD, and LRB policies. For these latter three, we used their authors’ C++ implementations available in [11], [7] and [40] respectively. We used their default configurations, with minor modifications, mainly to address trace parsing issues and remove redundant stats measuring. To verify a baseline, we make sure that LRU hit-ratio and byte-hit-ratio are the same in all three frameworks.

Traces Our measurements have used four types of real traces, taken from four different resources, as detailed in Table 1. MSR1, MSR2, and MSR3 are taken from an enterprise storage system [35]; SYSTOR1, SYSTOR2, and SYSTOR3 are taken from the storage of a VDI system [30]; CDN1 [9] along with CDN2 and CDN3 [41] are taken from CDN deployments; TENCENT1 [47] is taken from a large-scale photo service. Table 1 details for each trace the total number of accesses in the trace, the number of unique objects, and the total size of all unique objects.

Another interesting aspect of these traces is their object size distributions. This is depicted in Figure 8. As can be observed, all traces include objects of various sizes, ranging from less than 1KB to 0.5MB in the case of both storage systems, and up to 0.5GB or more in the case of the CDN. We can also see that in MSR1 and MSR2, it is easy to divide the objects into a small number of “size buckets”, suggesting that here slabbing might work well. In particular, 3–4 slabs could suffice. MSR3, TENCENT1 and SYSTORs traces are already more challenging, and the sizes of the objects are much less clustered. Finally, in the CDNs, object sizes are spread in the entire range of up to 0.5GB, hinting that a true size aware policy is needed.

5.1 Finding the Best Hit-Ratio Filter Variant

Figure 9 compares our AV with the IV approach (taken by Caffeine [32]) and the QV approach (used by Ristretto [18]). As illustrated, AV is consistently better for all traces in the entire range. Such improvement

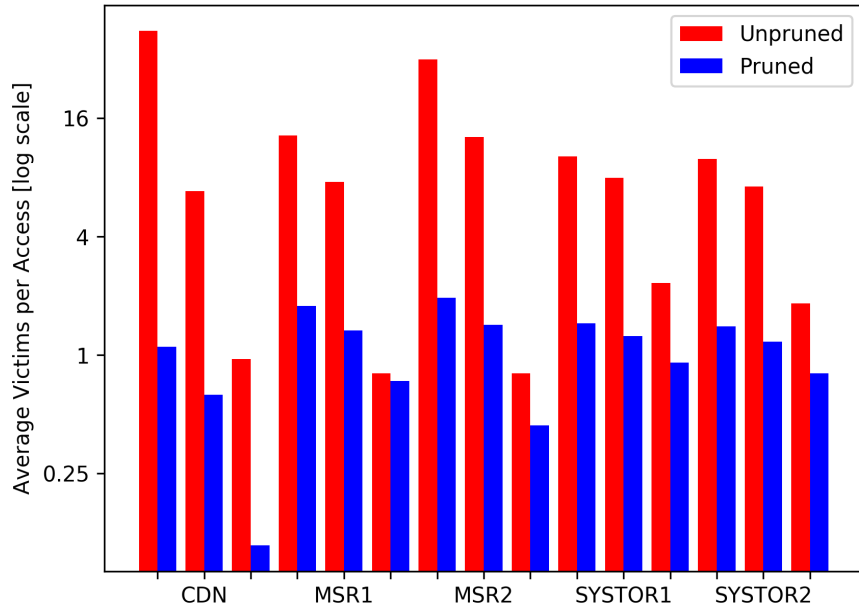


Figure 7: The number of victims per access, with and without the pruning optimization (see algorithm 4 line 7), for different traces with cache sizes of $10MB$, $1GB$, and $100GB$ for each.

is explained by the fact that AV more accurately captures the victim set’s benefit compared with IV and QV. Also, QV seems slightly better than IV.

Figure 10 repeats the above evaluation for byte hit ratio. Here, QV seems to be the best admission policy. Intuitively, AV’s relative weakness on this metric is because it is more likely to admit smaller items than larger items. Specifically, smaller items are likely to be compared against few victims, while larger items are compared against many victims. The byte hit ratio benefits more from serving a large item from the cache than from serving a small item. The QV approach is indifferent to items’ size and actively removes infrequent items from the cache even if it rejects the new item. These evictions free poorly utilized space, making QV more likely to admit large items to improve the byte hit ratio.

Notice that although there are clear differences between the admission approaches, the eviction policies make inconsistent difference, except for the Random policy that clearly lags. An interesting and expected observation is that Sampled Size is a little better than others for hit-rate due to the gain it gets from evicting large items, while it losses for byte-hit-rate. SLRU is decent for both measurements and more simple and efficient. Hence we continue with SLRU for all subsequent evaluations.

5.2 Hit-Ratios of State-of-the-Art Approaches

Next, we compare AV (SLRU) and QV (SLRU) with the seminal GDSF [13] as well as the recently proposed AdaptSize [10], and LHD [6], and LRB [41] schemes. Figure 11 shows the results for the hit ratio metric (important for user perceived reaction time). First, observe that LHD lags for small cache sizes, which is expected as LHD does not use any metadata on non-cached items. Intuitively, when the cache is large, LHD has enough information to reach sound decisions. But, this approach is insufficient for small cache sizes. In comparison, AV utilizes the TinyLFU data structure that collects metadata on a large sample of past

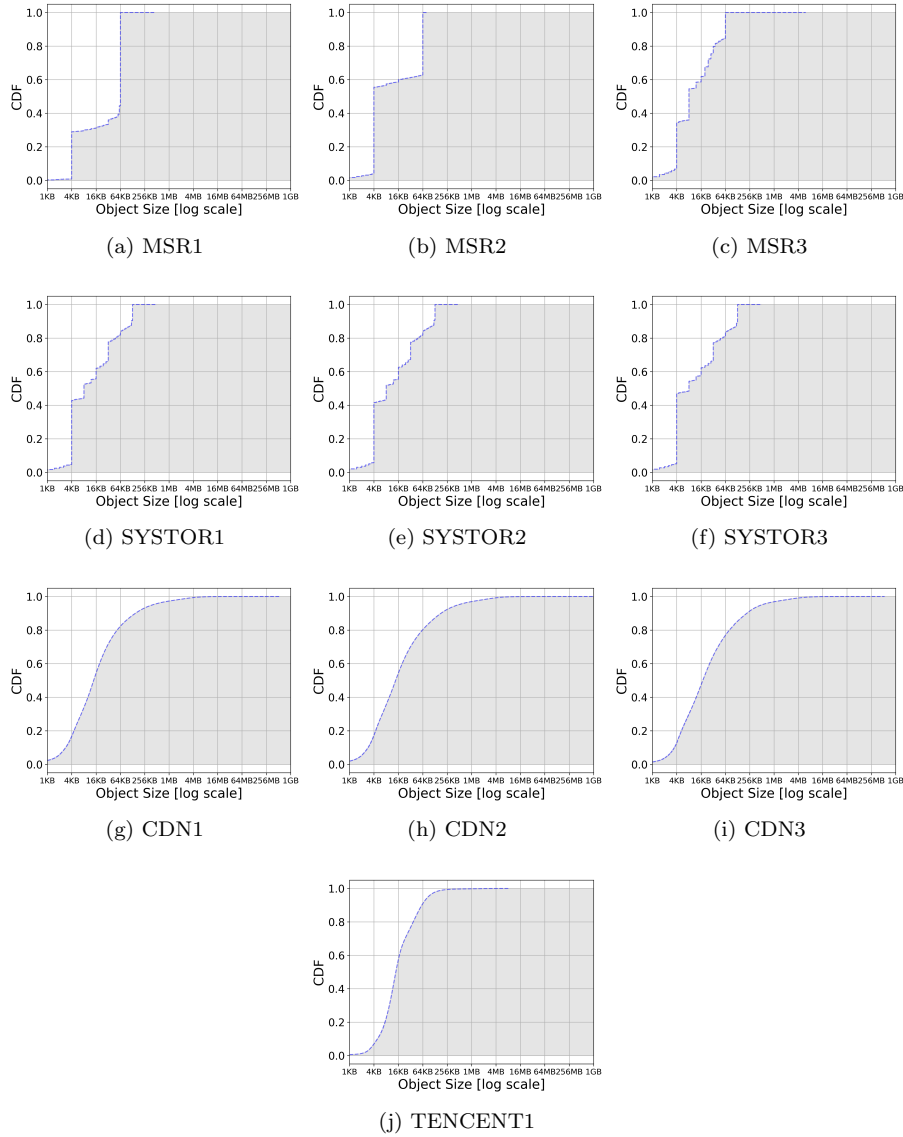


Figure 8: The cumulative distribution for objects sizes in each trace. The dotted blue lines indicate the actual span of objects sizes in the trace. Notice that in MSR1 and MSR2, most elements belong to a small number of tight object size ranges. In MSR3, TENCENT and the three SYSTOR traces, object sizes are spread further apart. Finally, in CDN traces objects sizes span the entire range.

accesses, which explains why AV outperforms LHD for small caches. We note that LRB is never the best policy despite its use of machine learning. This may not be surprising, as it is trained to maximize the byte hit ratio metric.

Next, observe that AdaptSize under-performs with large caches. Specifically, observe the MSR traces, the SYSTOR traces, and the CDN traces and notice that the hit ratio of AdaptSize does not improve when moving from 100GB to 10TB (an $\times 100$ increase). We believe that this is intrinsic to AdaptSize’s design. Specifically, AdaptSize admits a new item to the cache with a probability that is inversely proportional to the item’s size. AdaptSize is, therefore, unlikely to admit a very large item even if the cache has enough

Table 1: Summary of Traces

Name	Accesses <i>millions</i>	Objects <i>millions</i>	Total Objects Size
MSR1 [35]	29	18	738 <i>GB</i>
MSR2 [35]	37	6	161 <i>GB</i>
MSR3 [35]	2.2	0.27	5.3 <i>GB</i>
SYSTOR1 [30]	77	52	1.5 <i>TB</i>
SYSTOR2 [30]	79	48	1.4 <i>TB</i>
SYSTOR3 [30]	78	53	1.5 <i>TB</i>
CDN1 [9]	500	18	2.3 <i>TB</i>
CDN2 [41]	2,800	37	5.5 <i>TB</i>
CDN3 [41]	2,655	50	8.3 <i>TB</i>
TENCENT1 [47]	339	137	9.8 <i>TB</i>

Total number of accesses, unique objects and total size of unique objects for each trace.

space to accommodate this item without evictions. In that sense, AdaptSize is too selective and fails to utilize large caches.

We argue that AdaptSize can be improved for large caches if it bases the admission probability on the victim set’s size rather than on the cache candidate’s size. Such a change would bring it closer to our approaches in terms of design philosophy. That is, in AV the admission bar depends on the total frequency of the victim set. Similarly to AdaptSize, large items are less likely to enter the cache as they cause more evictions than small items do. Unlike AdaptSize, AV always admits an item if there is enough free space without evictions. Consequently, we can see that AV remains competitive across the entire range for all the traces.

The QV algorithm is not competitive in this metric since it evicts low-frequency items even when it fails to perform an admission. Thus, it has lower space utilization compared to AV. Finally, observe that GDSF remains competitive with present-day algorithms even after all these years. The main drawback of GDSF remains its logarithmic time complexity.

Figure 12 shows the results for the byte hit-ratio metric (important for bandwidth conservation). Notice that here, LRB is a very attractive suggestion, and is often one of the leading algorithms. AV is reasonable in the byte-hit-ratio metric as well but the QV algorithm excels in this metric. It never lags by much after the best alternative, and in some cases is strictly better. QV’s excellence in this metric is because QV actively evicts infrequent items even if the candidate is not admitted to the cache. Thus, it is more likely to pick up large and frequent items to increase the byte hit ratio.

5.3 Summary

As shown in our empirical study, AV is competitive with the best alternative (which is not always the same algorithm) in most workloads, especially for small to midsize caches. It is slightly better than the best alternatives in the byte hit-ratio metric.

We can further conclude that QV is inferior in the hit-ratio metric but almost always superior in terms byte hit-ratio. It is tough to excel in both hit ratio and byte hit ratio as there is some tension between them. For example, hit ratio wise it is better to cache an item of size 0.5MB that is accessed six times, but for byte hit ratio, it may be better to cache an object of size 1MB that is accessed only five times.

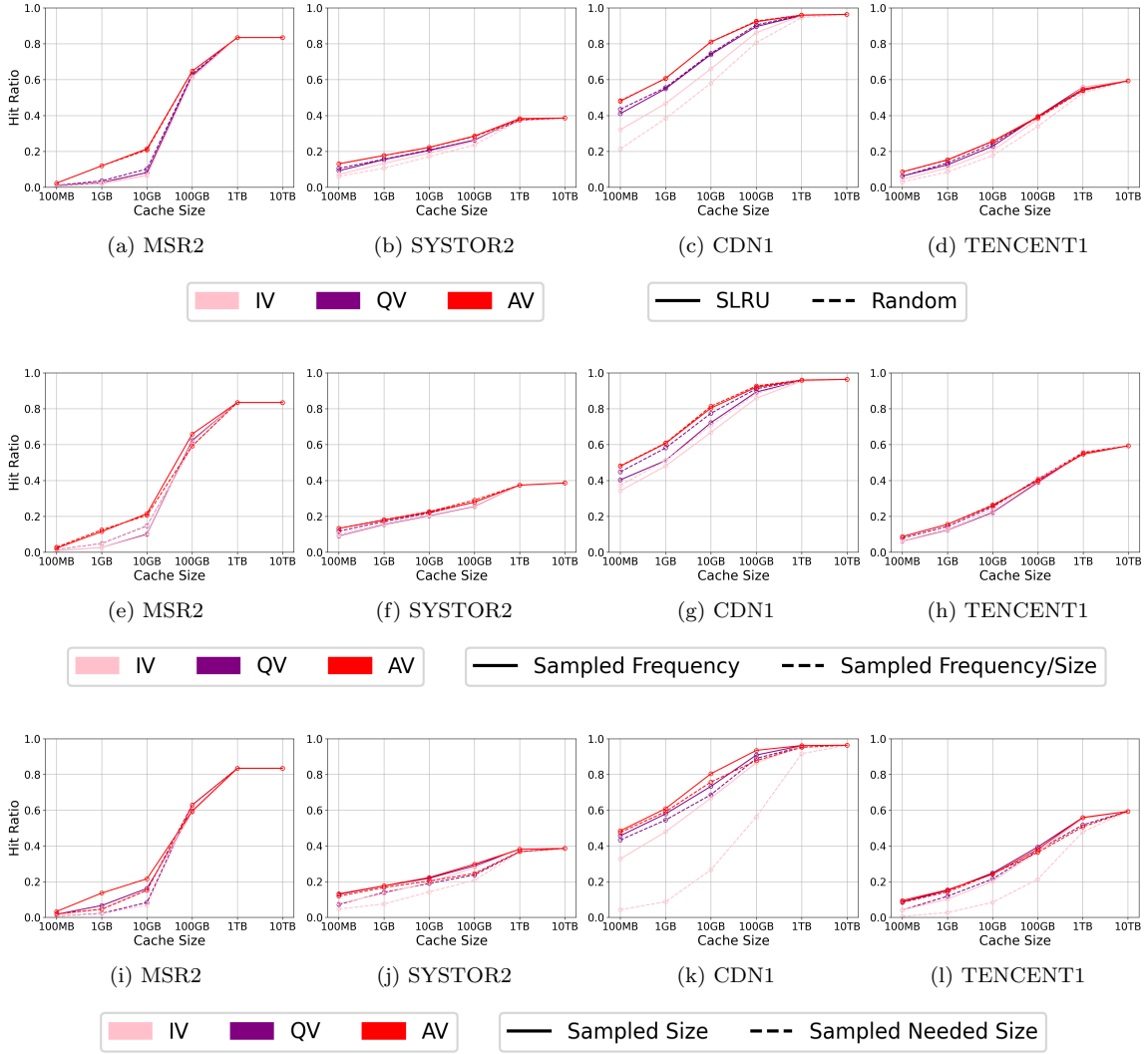


Figure 9: Evaluation of hit-ratio for W-TinyLFU versions. For brevity, we present here only MSR2, SYSTOR2, TENCENT1, and CDN1 traces. Results for the rest of the traces are similar.

5.4 Computational Performance

Finally, we complete the picture by measuring each policy’s algorithmic computation time when varying the cache size. Figure 13 and Table 2 illustrates these results. First, observe that AV (SLRU) is the most lightweight size-aware policy by a large margin and that its calculation time is comparable to the LRU policy.

One of the critiques about GDSF is that its logarithmic complexity in the number of cached objects makes it impractical for large caches. However, our evaluation shows that it is often considerably faster than both LHD and AdaptSize. We did notice the impact of the logarithmic complexity of GDSF, as its running time increased with the cache size, despite the higher hit-ratios obtained by large caches. This trend is evident in all traces, but is most acute for the CDN traces, which can be explained by their wide distribution of object sizes, as illustrated in Figure 8. Wider distributions of objects push the *actual average* running times closer to the logarithmic *upper bound* for GDSF’s priority queue maintenance. This is because when the size difference is very large, the distribution of objects’ priorities is larger, so objects are moved further

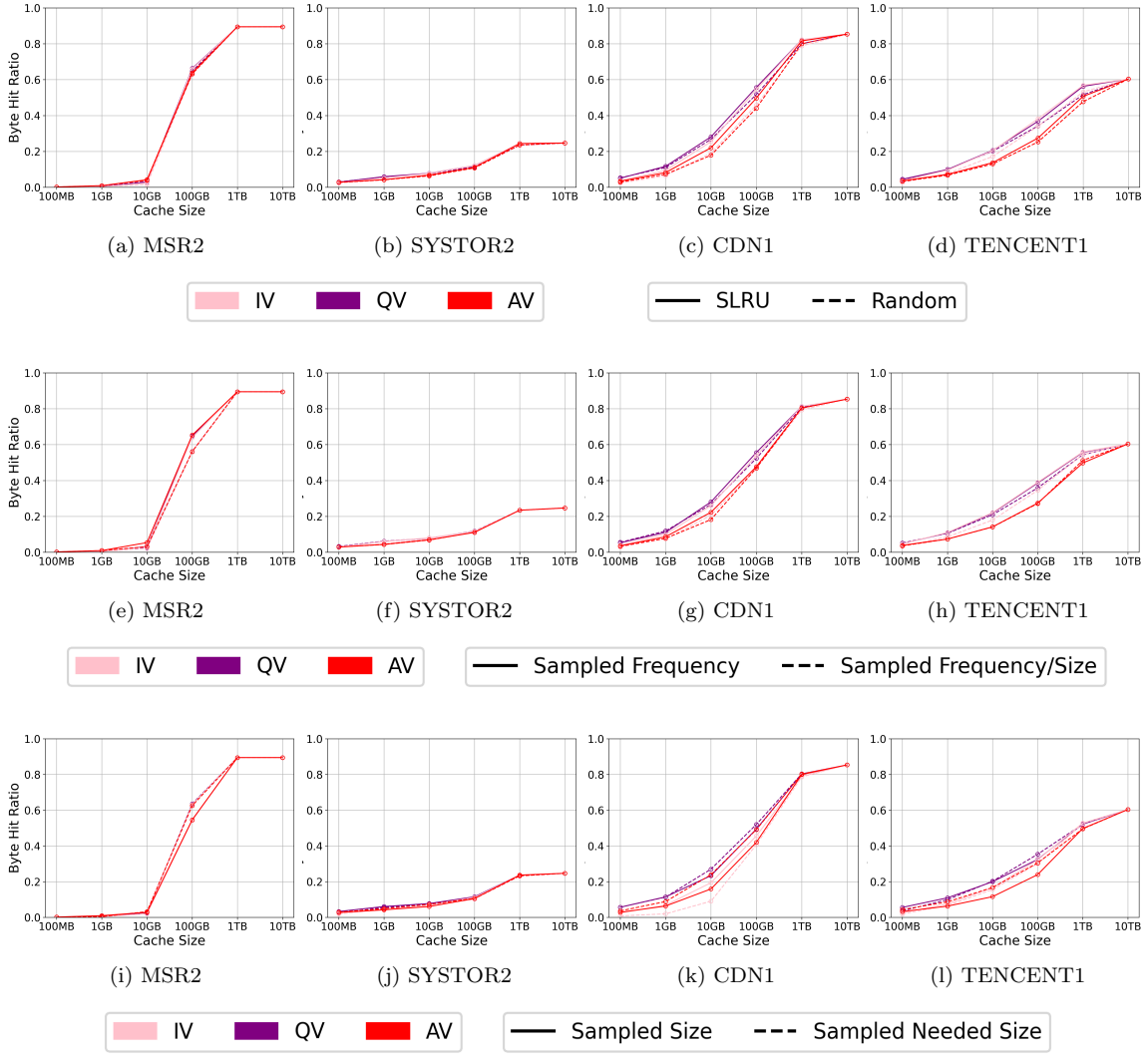


Figure 10: Evaluation of byte hit-ratio for W-TinyLFU versions. For brevity, we present here only MSR2, SYSTOR2, TENCENT1, and CDN1 traces. Results for the rest of the traces are similar.

inside the priority queue, which involves more internal operations.

Among the other algorithms, AdaptSize is the slowest one. Still, even LHD is at least x3 times slower than AV, especially where the hit-ratio is not very high, and many evictions are needed. Thus, we conclude that since neither of these algorithms yields a distinctive benefit in terms of hit-ratios, AV offers the most attractive trade-off as it is very simple to implement, as shown in Algorithms 1 and 4, and is considerably faster than the alternatives.

The run-time results for LRB (Table 2) are roughly an order of magnitude longer than the worst policy among the other schemes we tested (Figure 13). Yet, we notice that the performance improves when the cache size increases. In fact, LRB’s run-time is not directly dependent on the cache size. Rather, LRB works much harder on each miss than on a hit, as whenever there is a miss, it has to sample 64 items and invoke the machine learning based mechanism on them. Now, as the cache size increases, the hit ratio also improves, and therefore LRB does significantly less work. This is somewhat similar to LHD’s behavior, but

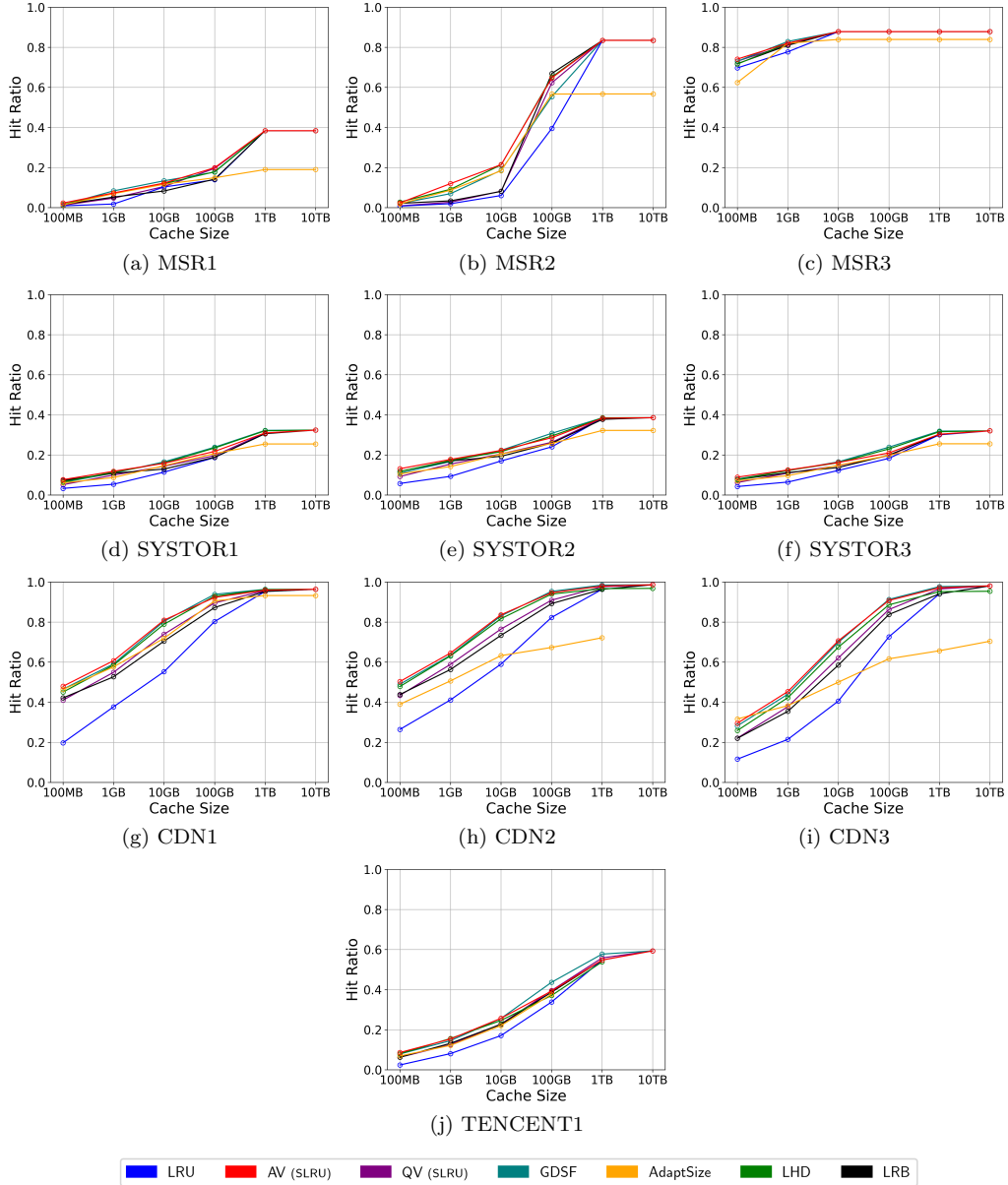


Figure 11: Evaluation of hit-ratio for state-of-the-art policies (given the traces’ characteristics in Table 1, a 1TB cache is practically an unbounded cache for the MSR traces and an almost unbounded for the SYSTOR traces; a 10TB cache is effectively an unbounded cache for all the traces).

on a grander scale.

6 Discussion

Modern databases and data-stores need to cope with objects of significant size variability. Therefore, traditional size oblivious cache management policies are often inadequate. In this paper, we have studied size aware cache management policies for workloads in which objects’ sizes vary. In particular, we focused on

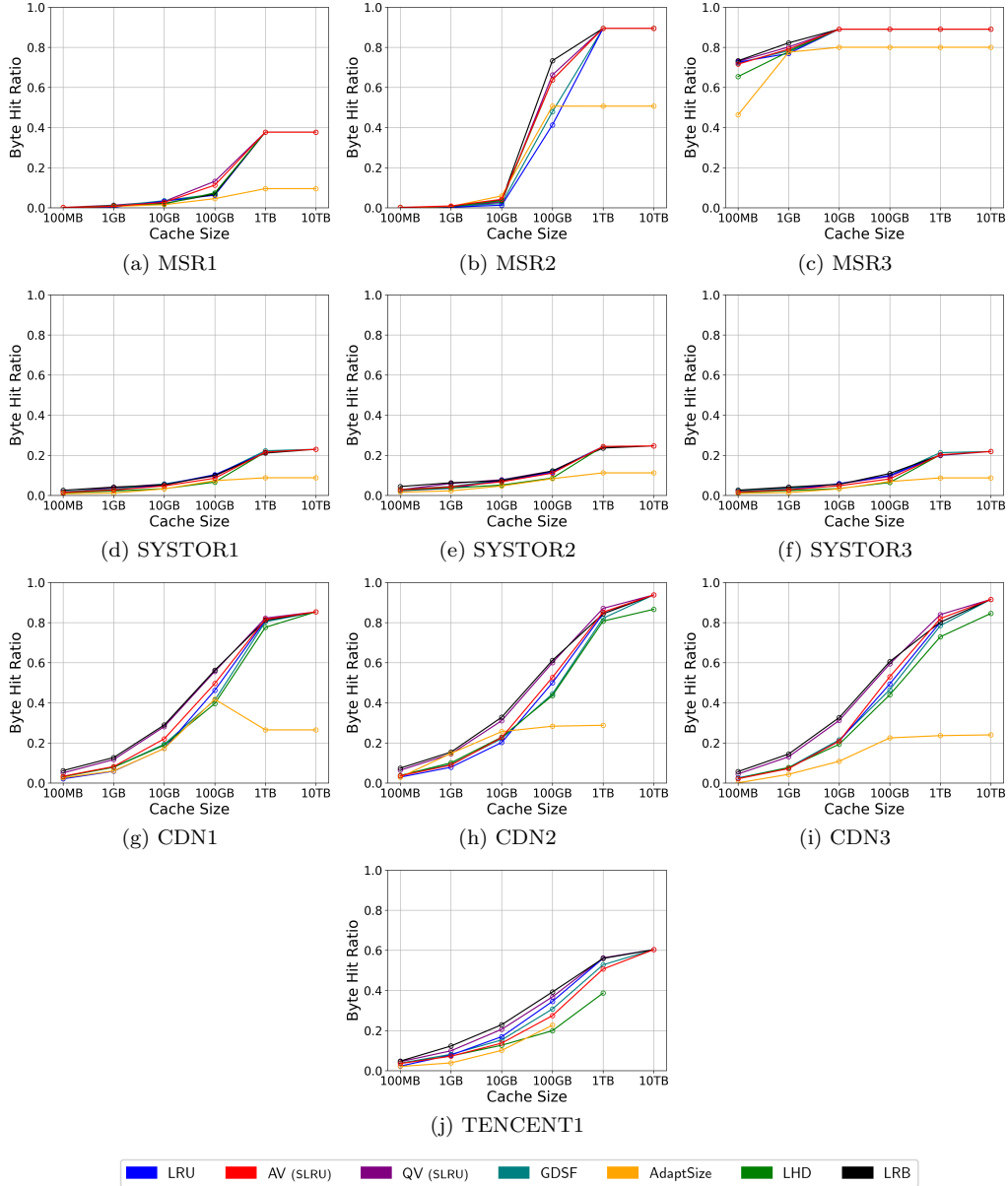
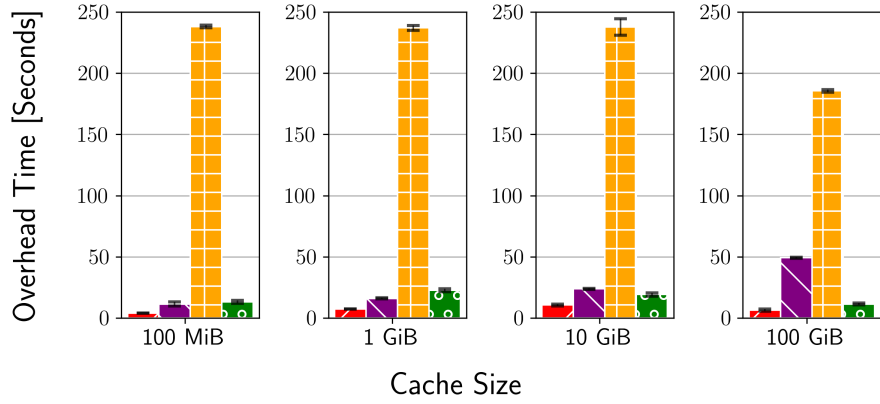
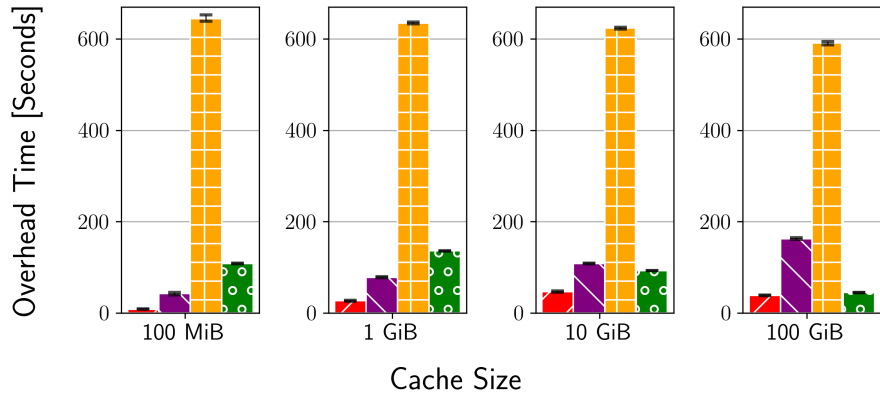


Figure 12: Evaluation of byte hit-ratio for state-of-the-art policies (given the traces’ characteristics in Table 1, a 1TB cache is practically an unbounded cache for the MSR traces and an almost unbounded for the SYSTOR traces; a 10TB cache is effectively an unbounded cache for all the traces).

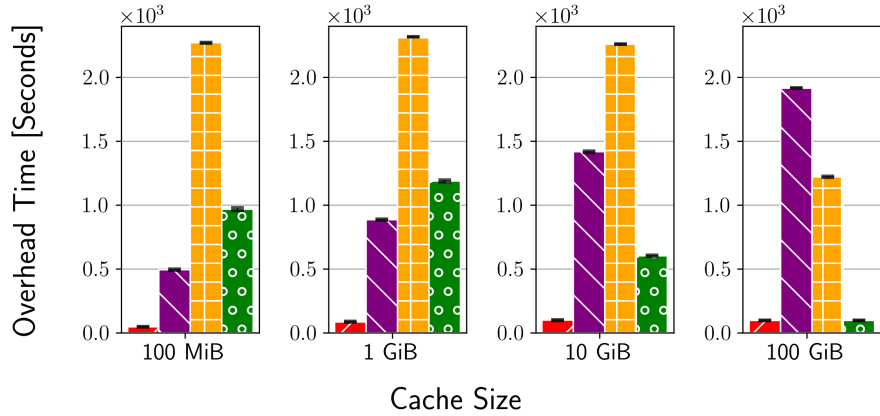
three approaches to extending the size oblivious W-TinyLFU policy into this more challenging setting. The first two, IV and QV, have been realized in two widely used open-source caching libraries [32, 18], but were never studied before, and the third is our new approach (AV). We also compared them against the state-of-the-art approaches GDSF [13], AdaptSize [10], LHD [6] and LRB [41]. We showed that in terms of hit-ratio and byte hit-ratio, none of these policies always wins. Yet, AV is always competitive with the state-of-the-art and is considerably faster than previous approaches. On the other hand, QV is less competitive on the hit-ratio metric but offers very attractive byte hit-ratio. In terms of CPU overhead, QV and AV are consistently the fastest.



(a) MSR1



(b) SYSTOR1



(c) CDN1



Figure 13: Overhead run-time of each policy for different traces and cache sizes. We ran GDSF and AV with our own Java implementation, while LHD and AdaptSize with their authors C++ implementation. To minimize the effect of using different frameworks and present the policy overhead only, we subtract the mean run-time of LRU in the matching implementation from the total run-time.

Trace	Cache Size	Overhead Time <i>seconds</i>
MSR1	100 MB	1,744
	1 GB	2,206
	10 GB	2,133
	100 GB	1,999
SYSTOR1	100 MB	9,666
	1 GB	8,257
	10 GB	7,410
	100 GB	6,380
CDN1	100 MB	26,744
	1 GB	25,725
	10 GB	16,712
	100 GB	14,769

Table 2: Overhead run-time of LRB with the authors’ C++ implementation. To minimize the effect of using different frameworks and present the policy overhead only, we subtract the mean run-time of LRU in the matching implementation from the total run-time.

Being fast is vital since caching is a performance optimization, so quicker policies are applicable in more systems and have a smaller CPU footprint. Thus, we believe that AV is a strong candidate as its hit-ratios are competitive, and it is considerably faster than the existing alternatives.

However, in scenarios where byte hit-ratio is more beneficial than hit-ratio (e.g., to minimize network bandwidth consumption), QV is likely the best alternative. Further, AV and QV are also very straightforward to implement. That is, their code is short, easy to explain, and requires no sophisticated math nor heavy machine learning mechanisms.

Our performance evaluation study has also revealed the following inherent drawback of AdaptSize [10], as a representative of policies that admit objects with a probability that is inversely proportional to the object’s size [1, 42, 36]. AdaptSize practically avoids inserting very large objects into the cache, even when it is far from full. Consequently, it fails to utilize significant portions of large caches, which hurts AdaptSize’s hit ratios, and even more significantly, its byte hit ratios when operating with very large caches.

Let us note that slabbing is orthogonal to the management policy. Yet, intuitively, the more effective the management policy is inside a slab, the fewer slabs one needs, and therefore one can expect to obtain a better overall hit ratio. Rigorously studying the tradeoffs of combining our policies with slabbing is left for future work.

Code availability The full code of the algorithms suggested in this paper is available at https://github.com/ohadeytan/caffeine/tree/arXiv_submission. The traces are available online and are pointed to from our code repository.

References

- [1] M. Abrams, C. R. Standridge, G. Abdulla, S. Williams, and E. A. Fox. Caching Proxies: Limitations and Potentials. Technical report, USA, 1995.
- [2] D. Adas, G. Einziger, and R. Friedman. TinyCache - An Effective Cache Admission Filter. In *GLOBECOM 2017 - 2017 IEEE Global Communications Conference*, pages 1–6, 2017.
- [3] Z. Akhtar, Y. Li, R. Govindan, E. Halepovic, S. Hao, Y. Liu, and S. Sen. AViC: A Cache for Adaptive Bitrate Video. In *Proc. of the 15th International ACM Conference on Emerging Networking Experiments And Technologies (CoNEXT)*, page 305–317, 2019.
- [4] M. Arlitt, R. Friedrich, and T. Jin. Performance Evaluation of Web Proxy Cache Replacement Policies. *Perform. Eval.*, 39(1-4):149–164, Feb. 2000.
- [5] S. Bansal and D. S. Modha. CAR: Clock with Adaptive Replacement. In *Proc. of the 3rd USENIX Conf. on File and Storage Technologies (FAST)*, pages 187–200, 2004.
- [6] N. Beckmann, H. Chen, and A. Cidon. LHD: Improving Cache Hit Rate by Maximizing Hit Density. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 389–403, 2018.
- [7] N. Beckmann, H. Chen, and A. Cidon. Simulation Code for the LHD Cache Replacement Policy. <https://github.com/CMU-CORGI/LHD>, 2018.
- [8] L. A. Belady. A Study of Replacement Algorithms for a Virtual-Storage Computer. *IBM Systems Journal*, 5(2):78–101, 1966.
- [9] D. S. Berger, N. Beckmann, and M. Harchol-Balter. Practical Bounds on Optimal Caching with Variable Object Sizes. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 2(2):1–38, 2018.
- [10] D. S. Berger, R. K. Sitaraman, and M. Harchol-Balter. AdaptSize: Orchestrating the Hot Object Memory Cache in a Content Delivery Network. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 483–498, 2017.
- [11] D. S. Berger, R. K. Sitaraman, and M. Harchol-Balter. Webcachesim, a C++ Cache Simulator. <https://github.com/dasebe/webcachesim>, 2017.
- [12] A. Blankstein, S. Sen, and M. J. Freedman. Hyperbolic Caching: Flexible Caching for Web Applications. In *2017 USENIX Annual Technical Conference (USENIX ATC)*, pages 499–511, 2017.
- [13] L. Cherkasova. Improving WWW Proxies Performance with Greedy-Dual-Size-Frequency Caching Policy. Technical report, In HP Tech. Report, 1998.
- [14] G. Chockler, G. Laden, and Y. Vigfusson. Data Caching as a Cloud Service. In *Proceedings of the 4th ACM International Workshop on Large Scale Distributed Systems and Middleware (LADIS)*, LADIS '10, pages 18–21. ACM, 2010.
- [15] G. Chockler, G. Laden, and Y. Vigfusson. Design and Implementation of Caching Services in the Cloud. *IBM Journal of Research and Development*, 55(6):9:1–9:11, 2011.
- [16] S. Cohen and Y. Matias. Spectral Bloom Filters. In *Proc. of the ACM SIGMOD Int. Conf. on Management of data*, pages 241–252, 2003.
- [17] G. Cormode and S. Muthukrishnan. An Improved Data Stream Summary: The Count-min Sketch and Its Applications. *J. Algorithms*, 55(1):58–75, Apr. 2005.

- [18] Dgraph. Ristretto: A High Performance Memory-Bound Go Cache. <https://github.com/dgraph-io/ristretto>, 2020.
- [19] G. Einziger, O. Eytan, R. Friedman, and B. Manes. Adaptive Software Cache Management. In *Proceedings of the 19th International Middleware Conference*, pages 94–106, 2018.
- [20] G. Einziger, R. Friedman, and B. Manes. TinyLFU: A Highly Efficient Cache Admission Policy. *ACM Transactions on Storage (TOS)*, 2017.
- [21] H. Gold. Netflix and YouTube are Slowing Down in Europe to Keep the Internet from Breaking. *CNN Business*, 2020.
- [22] S. Hasan, S. Gorinsky, C. Dovrolis, and R. K. Sitaraman. Trade-offs in Optimizing the Cache Deployments of CDNs. In *IEEE Conference on Computer Communications (INFOCOM)*, pages 460–468, 2014.
- [23] J. L. Hennessy and D. A. Patterson. *Computer Architecture - A Quantitative Approach (5. ed.)*. Morgan Kaufmann, 2012.
- [24] S. Jiang, F. Chen, and X. Zhang. CLOCK-Pro: an Effective Improvement of the CLOCK Replacement. In *Proc. of the USENIX Annual Technical Conference (ATC)*, 2005.
- [25] S. Jiang and X. Zhang. LIRS: an Efficient Low Inter-Reference Recency Set Replacement Policy to Improve Buffer Cache Performance. In *Proc. of the International Conference on Measurements and Modeling of Computer Systems SIGMETRICS*, pages 31–42, June 2002.
- [26] S. Jiang and X. Zhang. LIRS: An Efficient Low Inter-reference Recency Set Replacement Policy to Improve Buffer Cache Performance. *ACM SIGMETRICS*, pages 31–42, 2002.
- [27] T. Johnson and D. Shasha. 2Q: A Low Overhead High Performance Buffer Management Replacement Algorithm. In *Proc. of the 20th Int. Conf. on Very Large Data Bases (VLDB)*, pages 439–450, 1994.
- [28] G. Karakostas and D. N. Serpanos. Exploitation of Different Types of Locality for Web Caches. In *Proc. of the 7th Int. Symposium on Computers and Communications (ISCC)*, 2002.
- [29] A. Ketan Shah and M. D. Matani. An O(1) Algorithm for Implementing the LFU Cache Eviction Scheme. Technical report, 2010. "http://dhruvbird.com/lfu.pdf".
- [30] C. Lee, T. Kumano, T. Matsuki, H. Endo, N. Fukumoto, and M. Sugawara. Understanding Storage Traffic Characteristics on Enterprise Virtual Desktop Infrastructure. In *Proceedings of the 10th ACM International Systems and Storage Conference*, pages 1–11, 2017.
- [31] C. Li. DLIRS: Improving Low Inter-Reference Recency Set Cache Replacement Policy with Dynamics. In *Proc. of the 11th ACM International Systems and Storage Conference (SYSTOR)*, June 2018.
- [32] B. Manes. Caffeine: A High Performance Caching Library for Java 8. <https://github.com/ben-manes/caffeine>, 2016.
- [33] N. Megiddo and D. S. Modha. ARC: A Self-Tuning, Low Overhead Replacement Cache. In *Proc. of the 2nd USENIX Conf. on File and Storage Technologies (FAST)*, pages 115–130, 2003.
- [34] J. Mertz and I. Nunes. Understanding Application-Level Caching in Web Applications: A Comprehensive Introduction and Survey of State-of-the-Art Approaches. *ACM Comput. Surv.*, 50(6), Nov. 2017.
- [35] D. Narayanan, A. Donnelly, and A. Rowstron. Write Off-loading: Practical Power Management for Enterprise Storage. *ACM Transactions on Storage (TOS)*, 4(3):1–23, 2008.

- [36] G. Neglia, D. Carra, M. Feng, V. Janardhan, P. Michiardi, and D. Tsigkari. Access-Time-Aware Cache Algorithms. *ACM Trans. Model. Perform. Eval. Comput. Syst.*, 2(4), Nov. 2017.
- [37] E. J. O’Neil, P. E. O’Neil, and G. Weikum. The LRU-K Page Replacement Algorithm for Database Disk Buffering. *ACM SIGMOD Rec.*, 22(2):297–306, June 1993.
- [38] S. Park and C. Park. FRD: A Filtering Based Buffer Cache Algorithm that Considers both Frequency and Reuse Distance. In *Proc. of the 33rd IEEE International Conference on Massive Storage Systems and Technology (MSST)*, 2017.
- [39] Redis-Labs. Using Redis as an LRU cache. <https://redis.io/topics/lru-cache>, 2020.
- [40] Z. Song. A Simulator for CDN Caching and Web Caching Policies. <https://github.com/sunnyszy/lrb>, 2020.
- [41] Z. Song, D. S. Berger, K. Li, and W. Lloyd. Learning Relaxed Belady for Content Distribution Network Caching. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 529–544, 2020.
- [42] D. Starobinski and D. Tse. Probabilistic Methods for Web Caching. *Performance Evaluation*, 46(2):125 – 137, 2001.
- [43] R. Subramanian, Y. Smaragdakis, and G. H. Loh. Adaptive Caches: Effective Shaping of Cache Behavior to Workloads. In *Proc. of the 39th Annual IEEE/ACM International Symposium on Microarchitecture, MICRO*, pages 385–396, 2006.
- [44] A. S. Tanenbaum. *Modern Operating Systems (2nd ed.)*. Upper Saddle River, NJ, USA: Prentice-Hall, 2001.
- [45] C. Waldspurger, T. Saemundsson, I. Ahmad, and N. Park. Cache Modeling and Optimization using Miniature Simulations. In *USENIX Annual Technical Conference (ATC)*, pages 487–498, 2017.
- [46] N. Young. On-Line Caching as Cache Size Varies. In *Proc. of the second annual ACM-SIAM symposium on Discrete algorithms (SODA)*, pages 241–250, 1991.
- [47] K. Zhou, S. Sun, H. Wang, P. Huang, X. He, R. Lan, W. Li, W. Liu, and T. Yang. Demystifying Cache Policies for Photo Stores at Scale: A Tencent Case Study. In *Proceedings of the 2018 International Conference on Supercomputing*, pages 284–294, 2018.