

Adaptive Page Migration Policy With Huge Pages in Tiered Memory Systems

Taekyung Heo [✉], Yang Wang, *Student Member, IEEE*, Wei Cui, Jaehyuk Huh [✉], *Member, IEEE*, and Lintao Zhang, *Senior Member, IEEE*

Abstract—To accommodate the growing demand for memory capacity in a cost-effective way, multiple types of memory are incorporated in a single system. In such tiered memory systems consisting of small fast and large slow memory components, accurately identifying the performance importance of pages is critical to properly migrate hot pages to fast memory. Meanwhile, growing address translation cost due to the increasing memory footprints, helped adopting huge pages in common systems. Although such page hotness identification problems have existed for a long time, this article revisits the problem in the new context of tiered memory systems and huge pages. This article first investigates the memory locality behaviors of applications with three potential migration policies, least-recently-used (LRU), least-frequently-used (LFU), and random with huge pages. The evaluation shows that none of the three migration policies excel the others, as the effectiveness of each policy depends on application behaviors. In addition, the results show huge pages can be effective even with page migration, if a proper migration policy is used. Based on the observation, this paper proposes a novel dynamic policy selection mechanism, which identifies the best migration policy for a given workload. It allows multiple concurrently running workloads to adopt different policies. To find the optimal one for each workload, this study first identifies key features that must be inferred from limited approximate memory access information collected using accessed bits in page tables. In addition, it proposes a parallel emulation of alternative policies to assess the benefit of possible alternatives. The proposed dynamic policy selection can achieve 23.8percent performance improvement compared to a prior approximate mechanism based on LRU lists in Linux systems.

Index Terms—Tiered memory, page hotness, page migration, huge pages

1 INTRODUCTION

MEMORY systems are adopting multiple types of memory with different performance and capacity characteristics to increase the memory size in a cost-effective way. More expensive fast memory is backed by slower but higher capacity memory components, forming tiered memory systems. For capacity-optimized slow memory, non-volatile memories are already commercially available in the market, offering ~ 300 ns latency [1], [2]. In such tiered memory systems, accurately identifying the hotness of pages is critical to take advantage of the low-cost emerging memories while minimizing the performance loss.

Such page migration or replacement policies have been extensively studied for hardware caches, storage caches, and page management for virtual memory supports [3], [4], [5], [6], [7], [8], [9], [10], [11]. However, the emerging tiered memory systems provide new environments which are different from the prior work. The prior page replacement problems are restricted for choosing which data should be selected as victims to make room for new data. An access to

the slow component triggers an immediate promotion of the data to the fast component (cache). Unlike such replacement problems, the page migration problem must find which data must reside in the fast memory at a given time, while the data in the slow memory are accessible not necessarily triggering promotions.

In addition, compared to the storage caching systems, memory access traces collected from accessed bits in page tables are quite incomplete and approximate by their nature. Compared to the traditional page swapping, the tiered memory systems have much smaller latency and bandwidth differences between fast and slow memories than those of the prior DRAM and storage devices. Migrations between the two memory types are much more frequent than the prior studies as the costs of migration are relatively low. The new environments raise the need for revisiting the policy space of page migration in the context of tiered memory systems.

In the mean time, increasing memory footprints caused excessive misses in translation lookaside buffers (TLBs) for address translation. To reduce the cost of address translation, 2 MB huge pages have been adopted in x86 systems, which reduces TLB misses significantly for applications which otherwise suffer from the costs of TLB misses. While huge pages can drastically improve the address translation efficiency, their interaction with migration policies for tiered memory systems is yet to be investigated. Improving the accuracy of hotness measurement for huge pages and the efficiency of page migration for huge pages for tiered memory systems have been investigated by the prior studies [12],

• Taekyung Heo and Jaehyuk Huh are with KAIST, Daejeon 34141, South Korea. E-mail: {taekyung-heo, jhuh}@kaist.ac.kr.

• Yang Wang, Wei Cui, and Lintao Zhang are with Microsoft Research Asia, Beijing 100080, China. E-mail: {t-yangwa, weicu, lintaoz}@microsoft.com.

Manuscript received 20 May 2020; revised 7 Oct. 2020; accepted 1 Nov. 2020.

Date of publication 9 Nov. 2020; date of current version 13 Dec. 2021.

(Corresponding author: Jaehyuk Huh.)

Recommended for acceptance by A. R. Alameldoen.

Digital Object Identifier no. 10.1109/TC.2020.3036686

[13]. This paper explores the migration policy space with huge pages.

With such huge pages, this paper first investigates three widely accepted policies for page migration in tiered memory systems. It evaluates least-recently-used (LRU), least-frequently-used (LFU), and random policies for their behaviors on a range of applications. (i) Our investigation first shows that none of the policies excels the others for the range of applications, since each application has a different memory access behavior preferring a different migration policy. However, the page migration policies used by recent tiered memory studies are fixed to a single policy, such as a variant of LRU lists in Linux systems [13]. (ii) If a proper migration policy for each application is used, our investigation shows that huge pages can be effectively used for migration in tiered memory systems, reaping the benefit of efficient address translation.

Based on the observation of the migration policy evaluation, this paper proposes a new dynamic policy selection technique tuned for huge pages in tiered memory systems, called *Adaptive Migration Policy (AMP)*. AMP constantly collects memory access information using accessed bits in page tables, and periodically selects the best policy out of the three potential policies. To identify the best policy for a given workload, the study first identifies which features of memory access behaviors are highly correlated to the policy selection.

The first feature is used to identify workloads favoring random page placements. When accessed pages exceed the fast memory capacity, the locality cannot be captured effectively with the tiered memory system. Therefore, the random policy, which does not migrate pages actively, results in the best performance without migration overheads. If the workload exhibits a certain level of locality, either LRU or LFU is selected. To select the best one out of the two policies, AMP maintains a shadow page location states, which virtually mimics page migration. For example, if the current best policy is not LFU, the LFU policy is emulated with the shadow page location, and its effectiveness is tracked with the emulation. Based on the estimated effectiveness from the shadow state and the measured effectiveness from the current memory state, the better policy is selected for the next round.

Since the proposed technique can be applied for each memory control group (memcgs) in Linux, workloads in different memory control groups can select their own best policy. It allows the consolidated system to choose per-group optimal migration policies, allowing fine-tuning migration policies for co-running workloads.

We implement AMP in a Linux system, spanning from the kernel modification to the user-level components. The kernel is modified to track the recency and frequency of page accesses with accessed bits in page tables, and to provide policy options. The user-level components evaluate the features for each memory control group and apply the best policy for each group periodically.

We evaluate AMP in a Linux system with an emulated tiered memory system. The memory bandwidth of the slow memory node is throttled and saturated to emulate the slow memory. AMP can improve the performance of selected applications by 23.8 percent compared to the LRU lists

adopted in the prior work [13]. Furthermore, AMP can achieve 10.9, 6.4, 17.6 percent higher performance than LRU, LFU, and Random, respectively.

The contributions of our study are as follows:

- We find that workloads have diverse preferences on page migration policies in tiered memory systems, and we analyze the reason behind the page migration policy preferences.
- Our investigation shows that huge pages can be effective with page migration in tiered memory systems, if a proper migration policy is used.
- We define several features that have a relationship with the performance of page migration policies: fast memory hit ratio, page migration stability, and accessed page ratio. After that, we analyze the correlation between the features and the performance.
- We propose AMP, which dynamically selects a page migration policy between LRU, LFU, and Random using the features.

The rest of the paper is organized as follows. Section 2 describes the background of tiered memory and page migration policies. Section 3 investigates the behaviors of different page migration policies in the tiered memory system. Section 4 analyzes the critical features for determining the best policy, and Section 5 presents the implementation. Section 6 presents the experimental results, and Section 7 discusses the remaining issues. Section 8 concludes the paper.

2 BACKGROUND AND RELATED WORK

This section discusses the page migration problem in tiered memory systems, with its differences from the prior cache replacement along with the adoption of huge pages.

2.1 Tiered Memory Systems and Huge Pages

Tiered Memory Systems. A memory system composed of memories with various performance characteristics is called a *tiered memory system* [13]. The future memory systems are expected to be tiered memory systems due to the scaling limit of DRAMs. To increase the capacity of memory systems, memory systems are adopting non-volatile memories [1], [14], memory disaggregation [15], [16], [17], [18], and memory compression [19]. Usually, a tiered memory system is composed of fast memory and slow memory. Fast memory has a shorter latency and higher bandwidth compared to slow memory. A tiered memory system can be managed in hardware or in software. In this study, we assume a tiered memory system where an OS is responsible for managing data between two memory types. Data can be migrated at the page granularity, and the OS makes the page location and migration decisions.

Huge Pages for Efficient Address Translation. Future memory systems are expected to have TBs of memory with various latency and bandwidth [20]. In a memory system with a huge amount of memory capacity, address translations become a critical problem. Modern computer systems adopt virtual memory for efficient memory management. Therefore, virtual addresses should be translated to physical addresses to access data, and the mappings are maintained in a page table. The page table resides in main memory, and

translation lookaside buffers (TLBs) cache page table entries to avoid costly memory accesses. The memory capacity that can be translated by a TLB is called *TLB reach*. The problem is that the TLB size is limited to thousands of entries to shorten the access latency to TLBs. This limits the TLB reach of TLBs. The TLB reach can be increased with huge pages [21], [22]. While a single TLB entry can cover a 4 KB address space with 4 KB base pages, a TLB entry can cover several MBs or GBs with huge pages.

Previously, huge pages had several performance issues such as the increased page fault latency, memory bloating, unfair huge page allocations, and losing the page sharing opportunities [23]. Thanks to the recent studies to mitigate the problems [23], [24], [25], huge pages are becoming a viable option. Moreover, a multi-threaded page migration mechanism for transparent huge pages [13] makes migrating pages at the huge page granularity feasible. Therefore, in this study, we use 2 MB huge pages as a default page migration unit.

2.2 Page Migration Policies in Tiered Memory Systems

In this paper, we define a *page migration policy* as a policy that decides page locations in tiered memory systems, and it is used interchangeably with page hotness selection. Page migration policies try to fill fast memory with performance-critical hot pages. The page migration problem differs from the traditional cache replacement or page swapping, since any access to the slow memory does not necessarily trigger the promotion of the page to the fast memory. In hardware caches, memory blocks are immediately inserted to the cache for handling misses. The page swapping also requires to move a swapped out page from the storage to the memory to resolve the page fault. Unlike the cache replacement problems, the page migration problem needs to address which part of memory should reside in the fast memory, while the data in the slow memory are still accessible without migrating to the fast memory.

In addition, the cost and performance characteristics of tiered memory systems should be considered in the design and implementation of page migration policies. For example, maintaining an LRU stack is costly for virtual memory and architectural caches. In virtual memory, CLOCK approximates LRU with a single reference bit [26]. In CPU caches, tree-based pseudo-LRU is used to lower the area overhead [27], [28]. Likewise, cache replacement policies should be adopted to tiered memory systems with the consideration of the cost and performance characteristics. In this section, we summarize the prior studies on page migration policies.

- ① *Recency-based policies.* Native Linux systems have LRU lists to reclaim pages under memory pressure. The native LRU lists approximate LRU with two LRU lists. One is the active list, and the other is the inactive list. The membership of pages is controlled by the heuristic implemented in the kernel. The kernel uses the accessed bits of pages and page types to update the membership of pages. Each page has an accessed bit in the page table entry (PTE). When a page is accessed, the corresponding accessed bit is

set. As the goal of LRU lists is to reclaim pages under memory shortage, the native LRU lists do not update the membership of pages in non-memory pressure conditions.

- ② *Frequency-based policies.* Access frequency has been widely adopted in the page management for tiered memory systems. In modern processors, the exact access frequency of a page cannot be obtained due to the lack of hardware support. Instead, access frequency can be estimated using accessed bits. The accessed bit of a page is periodically checked and recorded to a per-page bit vector. An access frequency can be calculated by averaging the number of bits set in the bit vector. HeteroVisor [29] tracks the access frequency of pages to decide page locations between fast die-stacked DRAMs and slow off-chip DRAMs. If the access frequency of a page exceeds a predefined hot page threshold, the page is classified as hot and migrated to the fast die-stacked DRAM. On the other hand, Thermostat [12] finds the access frequency threshold using a user-specified maximum allowable slowdown.
- ③ *Limiting the promotion rate.* In a tiered memory system where pages have to be migrated from slow memory to fast memory on every page access to slow memory, limiting the number of page migrations from slow memory to fast memory is a way to guarantee the performance slowdown [19]. The rate of pages migrated from slow memory to fast memory within a time window is defined as the *promotion rate*. Lagar-Cavilla *et al.* find the relationship between the page access recency and promotion rate, and the pages that are expected to show a low promotion rate are identified by measuring the access recency of a page.

2.3 Other Prior Work for Cache Replacement

In the prior cache replacement, unlike page migration, the promotion to the cache is immediately triggered while handling misses. Therefore, the hotness selection within the cache is focused on which data should be evicted as victims to make room for newly inserted data. For selecting the victims, there have been many different approaches with their own advantages and disadvantages. In this section, we summarize the prior studies on cache replacement policies.

- ① *Recency-based policies.* The most common cache replacement policy is the Least Recently Used (LRU) replacement policy [3]. LRU maintains an LRU stack, which is sorted by the access recency. The most recently accessed page is placed at the top of the LRU stack, and the least recently accessed page is placed at the bottom. Once a page is referenced, the page is removed from its position and moved to the top of the stack. On cache evictions, LRU replaces the least recently accessed page. Although LRU works well for most memory access patterns, LRU fails with scanning memory access patterns and loops. Scanning memory access patterns access pages only once and never use them again. Therefore, caching recently accessed pages wastes the cache space. Loops that access pages

larger than the cache size evict pages that are accessed in the near future. Prior studies try to mitigate the limitations of LRU by identifying the anomalies [30], [31], [32].

- ② *Frequency-based policies.* Least Frequently Used (LFU) tracks the access frequencies of pages and evicts a page with the least access frequency. Although LFU has an advantage that it can exploit the long-term access history, it has several disadvantages. LFU can make a wrong decision in the initial stage due to the lack of access history. Moreover, LFU may hold stale pages that are not hot anymore due to their previous access history.
- ③ *Adaptive policies.* There have been many prior studies to improve the hit ratio of caches [4], [7], [10], [33], [34], [35]. Among the prior studies, we focus on the adaptive selection of cache replacement policies [5], [6], [8], [9], [11], [36], [37]. Abstractly, ARC [8] and CAR [9] partition a cache into two partitions for recently accessed pages and frequently accessed pages. The partition sizes are dynamically adjusted based on the workloads' behavior. A few prior studies adopt machine learning techniques to choose a cache replacement policy adaptively [6], [37]. In CPU caches, set dueling has been proposed to dynamically choose a policy between two competing policies [11]. Cache sets are divided into dedicated sets and follower sets. A small number of sets are allocated to dedicated sets, and the dedicated sets are managed with two competing policies to evaluate the performance of the policies. The policy that performs better is applied to the follower sets.

3 MIGRATION POLICIES WITH HUGE PAGES

3.1 Migration Policies

To determine which pages should be placed in the fast memory, the page hotness selection mechanism consists of two components. First, it must record the access history of each huge page. The second component chooses hot pages based on the access histories of all pages. In this section, we investigate the page migration policies using huge pages (2 MB). Improving the accuracy of hotness measurement for huge pages and facilitating the page migration for huge pages for tiered memory systems have been investigated by the prior studies [12], [13]. This paper explores the migration policy space with huge pages.

Although the main reason for using huge pages is to mitigate the cost of TLB misses, it also allows reducing the complexity of page managements for migration. We track the access information for each huge page to reduce the overhead of tracking the information. In addition, it also reduces the latency to select hot pages, as the number of candidate pages is significantly reduced by huge pages. A possible downside of using huge pages in tiered memory system is the potential waste of some fast memory, when only a small subset of huge pages are actively accessed. However, this paper is focused on page management based on huge pages, since huge pages were effective across all benchmark applications in our study, as shown in Section 3.3.

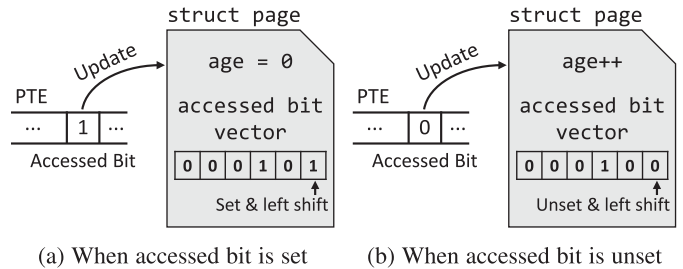


Fig. 1. Tracking page age and access frequency using accessed bits.

Tracking Access History. To record the access history of each page, we use the accessed bit in each page table entry, which is set by the processor hardware. Fig. 1 illustrates how the page history information of each page tracked. For each huge page, we record two types of access information, *age* and *history*. Every five seconds, the kernel checks the accessed bit of each huge page to update its age and history. Once the information is updated, the accessed bit in the page table entry is reset.

Age represents the recency of access for a page. As shown in the figure, if the accessed bit is set during the 5-second period, the page *age* is reset to 0. If the accessed bit is not set (no access for the last five seconds), its age is increased by 1. Therefore, the age of a page means when the last access for the page occurred at 5-second granularity.

History (accessed bit vector) represents the accessed bit vector for each page for the past n periods. Fig. 1 shows that six bits are used for the access history. A new accessed bit is pushed to the tail of the vector. Although the accessed bit vector can be used for representing *age*, we use a separate age variable to record the age of longer periods without incurring a long bit vector for each page. The detailed accessed bit vector covers a shorter time range than the age variable. In this paper, we use 64 bits for each huge page for the accessed bit vector.

As will be discussed later, the proposed memory manager allows adjusting how much fast memory and slow memory can be used for each application group. If the number of application groups is set to one, all applications share the entire fast and slow memory.

Every five seconds, the access information of all huge pages are updated, and the set of pages which must be in the fast memory are determined. For the pages which are selected for the fast memory, but not already in the fast memory, the migration of the pages are initiated. Since we use huge pages, this step does not incur significant overheads. First, the number of huge pages for a given application footprint is 512 times smaller than that of base pages. Therefore, tracking and sorting huge pages cause only a small extra overhead during the five second period. Second, if a good migration policy is selected, only a small fraction of pages are migrated to the fast memory, since it is likely that many hot pages are already in the fast memory.

Migration Policies. Among various page migration policies, we investigated the workloads' preference on policies of LRU, LFU, and Random. ① LRU tracks the access recency of pages with ages, and the most recently accessed pages are migrated to the fast memory. Using the age information for all pages, pages with the lowest ages are considered hot. ② LFU tracks the access frequency of pages, and the most

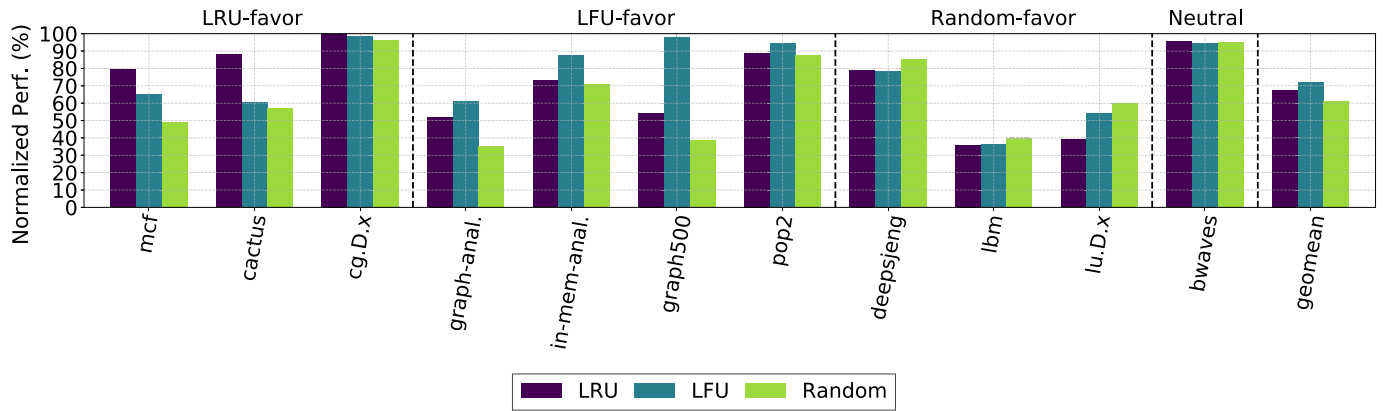


Fig. 2. Normalized performance of workloads with various page migration policies.

frequently accessed pages are located in fast memory. The access frequency is obtained with the per-page accessed bit vectors, by averaging the number of bits set in the bit vector.

③ **Random** simply fills the fast memory first, once the fast memory is filled, the slow memory is allocated. In Random, if a page located in the fast memory is freed, a random page from the slow memory is migrated to the fast memory.

3.2 Results

Methodology. We find the workloads' preferences on page migration policies by measuring the performance of workloads on an emulated tiered memory system. If a policy can accurately identify hot pages, the policy can present a higher performance compared to the other policies. We compose a tiered memory system on a NUMA system by throttling the memory bandwidth of one node. The bandwidth-throttled node becomes a slow memory node. 2 MB transparent huge pages are used, and pages are migrated at the huge page granularity. In the rest of paper, *fast memory ratio* is the ratio of the fast memory over the entire memory footprint, which is set to 50 percent in the evaluation of this section. 4B is used to track the page age, and 8B is used to track the access frequency. We evaluate the performance of page migration policies with selected workloads from SPEC CPU 2017 [38], CloudSuite [39], NPB [40], and graph500 [41]. In the remaining part of this section, we describe the reason behind the page migration policy preferences.

Fig. 2 presents the normalized performance of workloads with various page migration policies. We measure the execution time, and the performance is the reverse of execution time. The performance is normalized to the performance with the 100 percent fast memory ratio. Workloads can be classified into four groups: LRU-favor, LFU-favor, Random-favor, and neutral, as shown in the figure.

① **LRU-Favor.** Workloads with strided memory access patterns favor LRU. Strided memory access patterns sequentially access pages with the same distance between memory accesses with low data reuse. Therefore, keeping frequently accessed pages in fast memory may degrade the performance. LRU can keep the recently accessed pages in fast memory, increasing the probability of accessing fast memory. *mcf*, *cactus*, and *cg.D.x* have strided memory access patterns. *mcf* has a pointer chasing in `price_out_impl` as shown in Code 1 [42]. The pointer chasing pattern of *mcf* is actually a strided memory

access pattern because the data structures are sequentially located in a virtual address space [42], [43]. *cactus* strides over one dimension of a matrix while working on a multi-dimensional matrix [44]. *cg.D.x* calculates the eigenvalues of a sparse matrix using the conjugate gradient method. It operates on a large matrix, and the matrix is accessed sequentially with low data reuse [45].

Code 1. A Function That Shows a Strided Memory Access Pattern in *mcf*

```

1 long price_out_impl(network_t *net)
2 {
3   ...
4   iterator = first_list_elem->next;
5   while (iterator) {
6     arcin = iterator->arc;
7     tail = arcin->tail;
8     ...
9     iterator = iterator->next;
10  ... } ...
11  }
```

LFU-Favor. Workloads with frequently accessed data structures favor LFU. For the workloads, LFU can keep the frequently accessed hot data in fast memory, and it can protect the fast memory from being polluted by recently accessed cold data. *Graph-analytics* runs the PageRank algorithm, and it updates the ranks of neighbor vertices while walking on vertices. A vertex with more neighbors tends to be accessed more frequently. LFU can identify the vertices with more neighbors and keep them in fast memory. The first subfigure of Fig. 3 shows the temporal change of access frequencies of pages of *graph-analytics*. The access frequencies of pages are tracked by checking the accessed bits of pages on every one epoch, whose length is four seconds. The size of the per-page accessed bit vector is 8-bit. The pages with higher access frequencies are drawn at the bottom, and the pages with lower access frequencies are illustrated at the top of the figure. *Graph-analytics* has several scanning memory access patterns. We emphasize one of the scanning patterns with an ellipse in the figure. LRU fails to keep frequently accessed hot pages in fast memory due to the scanning patterns.

In-memory-analytics runs the alternating least squares algorithm. It trains a model multiple times with

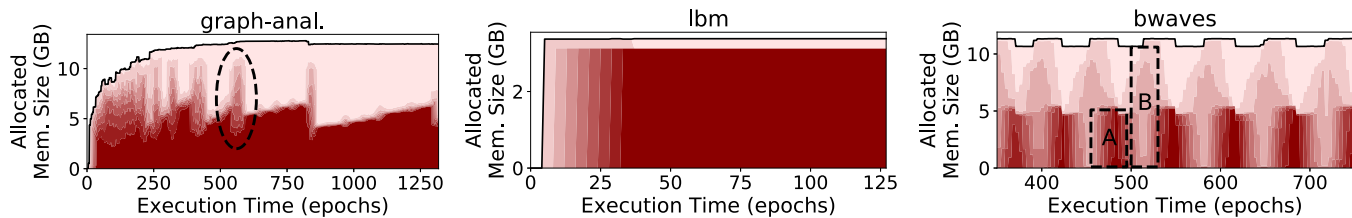


Fig. 3. Temporal change of access frequencies of pages. Pages are sorted by their access frequency, and the pages with higher access frequencies are drawn at the bottom with a darker color.

various parameters to find the best parameters. On every training, the training dataset is loaded by scanning the dataset. LFU can keep the frequently accessed trained model in fast memory. *graph500* runs the breadth-first search (BFS) algorithm multiple times. On every BFS, *graph500* validates the result with the `validate_result()` function. Therefore, the graph is accessed frequently, and *graph500* favors LFU.

Random-Favor. Workloads with low locality favor Random. First, workloads that have pages with similar access recency or frequency prefer Random. If all pages are equally recently accessed or frequently accessed, migrating pages with LRU or LFU adds the performance overhead without benefit, and Random can eliminate the overhead. *lbm* and *lu.D.x* prefer Random because of this reason. *lbm* runs the Lattice Boltzmann Method to simulate fluids. *lbm* allocates grids that represent three dimensions. *lbm* sweeps the grids multiple times within a short time to simulate fluid collisions. As a result, the pages of *lbm* show similar access recency and frequency. The second subfigure of Fig. 3 illustrates the temporal change of access frequencies of *lbm*, showing homogeneous page access frequencies. Second, workloads with random memory access patterns favor Random. LRU nor LFU cannot find hot pages from workloads with random memory access patterns. *deepsjeng* has a random memory access pattern. *deepsjeng* is a chess solver, and it has a hash table for the alpha-beta tree searching to find the next move. Memory accesses to the hash table show a random pattern.

Neutral. Workloads with mixed memory access patterns show a neutral preference on page migration policies. *bwaves* has mixed memory access patterns. The last subfigure of Fig. 3 shows the temporal change of access frequencies of *bwaves*. The execution of *bwaves* is composed of two phases. The first phase accesses the frequently accessed data continuously, favoring LFU (pattern A). The second

phase accesses the remaining working set with a strided memory access pattern, favoring LRU (pattern B). As a result, *bwaves* does not have a large performance gap between policies.

3.3 Huge Pages With the Prior Modified LRU Lists

A prior study proposed to reuse the LRU lists in the native Linux kernel to identify hot pages and migrate them to fast memory [13]. The study periodically scans the LRU lists to update the membership of pages. In the following sections, we call the LRU lists *the modified LRU lists*. The modified LRU lists classify pages in the active list as hot, and pages in the inactive list as cold. The modified LRU lists migrate the pages in the active list to fast memory and pages in the inactive list to slow memory. We use the modified LRU lists as the baseline migration policy to compare against the proposed technique in this paper.

In this subsection, we present that the performance of the modified LRU lists is parameter-sensitive. The modified LRU lists have two parameters: active list scanning ratio and inactive list scanning ratio. The (in)active list scanning ratio determines how many pages are scanned from the list on every scan. While scanning pages, the accessed bits of pages are checked, and the membership of pages is updated. If half of the pages in a list are scanned, the scanning ratio is 50 percent. By default, the modified LRU lists scan 50 percent of the active list and inactive list on every five seconds [46]. We evaluate the performance of workloads while varying the scanning ratios of lists from 0 to 100 percent with a 20 percent step. Experiments are run with 2 MB transparent huge pages.

Fig. 4 shows the normalized performance of selected workloads with various combinations of the active list scanning ratio and inactive list scanning ratio. We use the same experiment environment and definition for the performance that we use in Section 3.2. The cells with high performance

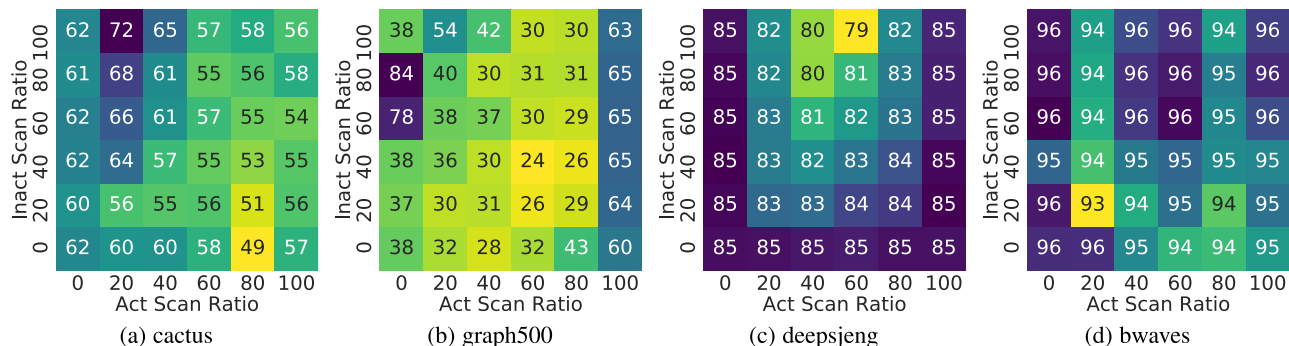


Fig. 4. Normalized performance of workloads with various parameter combinations when the modified LRU lists are used.

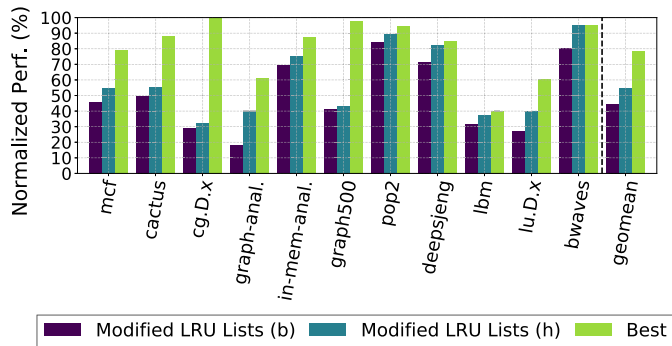


Fig. 5. Performance of the modified LRU lists with 4KB (base) and 2MB (huge) page sizes, compared to the potential best policy.

are drawn with a darker color, and the cells with low performance are drawn with a lighter color. The parameter combination that shows the best performance is different for each workload. Note that the performance gap between the best and worst parameter combinations is 60 percent in *graph500*, showing the parameter-sensitivity of the modified LRU lists.

The problem originates from the inappropriate application of LRU lists to a tiered memory system. The original goal of LRU lists is to reclaim pages under memory pressure. Therefore, the pages in the inactive list become page reclamation *candidates*. It does not mean that the pages in the inactive list are cold. In the native LRU lists, pages in the inactive list are moved back to the active list again if the pages are accessed. However, in the modified LRU lists, pages in the inactive list are migrated to slow memory, assuming that all pages in the inactive list are cold. Although the cold-classified pages will be migrated back to fast memory when they are accessed, the performance degradation cannot be avoided. Therefore, the latest proposal, the modified LRU lists are not the best option for tiered memory systems. In the following sections, we run the modified LRU lists with the default parameters (50%, 50%) [46].

Fig. 5 presents the performance of the modified LRU lists with the base page size (4 KB), and huge page size (2 MB). The performance is normalized to that with the ideal memory which consists of only fast memory. In addition to the performance of the modified LRU lists, it also shows the potential performance when the best replacement policy is selected among the aforementioned three policies. The results show that the huge page can improve the performance by 23.5 percent on average with the modified LRU lists. There are two potential advantages of huge pages. First, it reduces TLB misses by the increased translation capability. Second, page migration also uses huge pages,

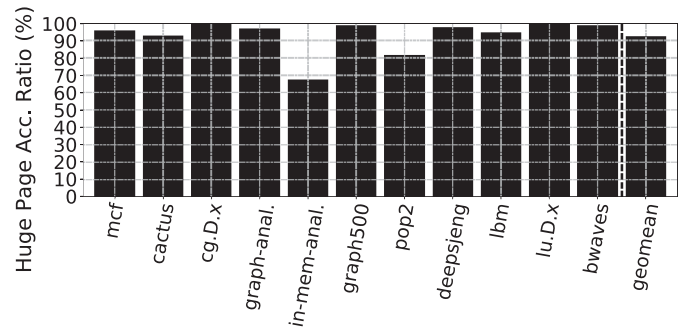


Fig. 7. Huge page access ratio of workloads.

and thus when spatial locality exists, it can prefetch a large chunk of data from the slow to fast memory. However, the modified LRU lists do not provide good identification of page hotness, compared to the potential best policy. The best policy can excel the modified LRU lists by 33.71 and 23.5 percent, compared to the modified LRU list with the base and huge page sizes, respectively.

3.4 The Homogeneity of Huge Page Hotness

One of the key requirements of migrating pages at the huge page granularity is the homogeneity of hotness in a huge page. To show how much of a 2 MB page is actually accessed, we present the results of the homogeneity of page hotness by measuring the number of accessed 4 KB base pages within an accessed huge page. We define the ratio as the huge page access ratio. If a workload exhibits a high level of hotness homogeneity, the rest of the base pages within a huge page are likely to be accessed, when a base page in the huge page is accessed. The time interval is determined by the multiplication of the accessed bit check interval (4-second) and the length of the per-page bit vector (8-bit). As the accessed bits of base pages in a huge page cannot be tracked in the current system, for this analysis, we use 4 KB base pages to track the accessed bits of pages. Based on the access statistics on base pages, we infer how many base pages within a huge are accessed in each time interval. Fig. 7 presents the huge page access ratio of workloads. On average, our workloads have huge page access ratios higher than 92 percent, justifying the huge-page-granular migrations.

4 ADAPTIVE PAGE MIGRATION POLICY SELECTOR

We present AMP, which adaptively selects a page migration policy preferred by a workload. AMP chooses a page migration policy between LRU, LFU, and Random. In this section,

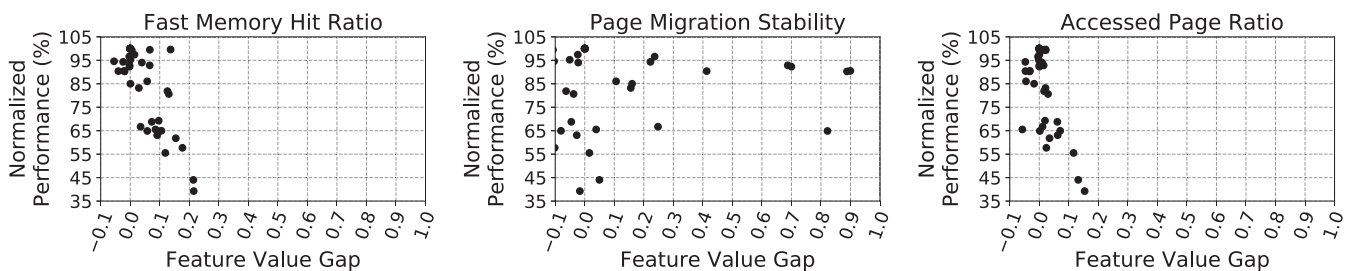


Fig. 6. Feature-performance scatter charts.

TABLE 1
Pearson Correlation Coefficients Between Feature Value Gap and Normalized Performance

	Correlation Coeff.	P-value
Fast Memory Hit Ratio	-0.8194	1.0505×10^{-11}
Migration Stability	0.0322	0.8355
Accessed Page Ratio	-0.4975	0.0006

we define features and analyze the relationship between features and performance. At last, we present the main algorithm of AMP based on the feature analysis.

4.1 Feature Analysis

Correlation Analysis. We define three features that are possibly related to the workloads' preferences on page migration policies: fast memory hit ratio, page migration stability, and accessed page ratio.

The fast memory hit ratio is the number of accessed pages in fast memory divided by the total number of pages. We assume that a page migration policy that can identify hot pages better shows a higher fast memory hit ratio.

The page migration stability is the number of stable pages divided by the number of total pages. The locations of pages are updated periodically in the baseline page migration policies. A page is regarded as *stable* if the page location has not been changed compared to the previous location. The page migration stability presents the page migration cost of a page migration policy.

The accessed page ratio is the number of accessed pages divided by the total number of pages. The assumption behind this feature is that workloads touch pages as they progress. If a page migration policy has been successful in choosing hot pages, a workload can progress faster. Therefore, the page migration policy may present a higher accessed page ratio. These features are measured on every page migration, whose interval is five seconds.

We analyze the relationship between features and the normalized performance of workloads using the Pearson correlation coefficient. We define a *feature value gap*, which is the gap between the feature value of the page migration policy that performs the best and the feature value of a selected page migration policy. If the feature plays an important role in the performance, the lower the gap is, the closer the performance of the selected policy is to the performance of the best-performing policy. The performance is the reverse of the execution time, and it is normalized to the best-performing page migration policy. We calculate the Pearson correlation coefficient between the feature value gap and the normalized performance. The analysis is conducted on the data that we present in Section 3.1.

Fig. 6 shows the scatter charts between the feature value gap and normalized performance. Table 1 presents the Pearson correlation coefficients. Additionally, it shows the p-values of the correlation coefficients. The absolute value of a correlation coefficient presents the strength of the correlation between the feature and performance. P-values show statistical significance. A feature is considered to have a statistically significant correlation if its p-value is lower than 0.01. Among the features that we have defined, the fast

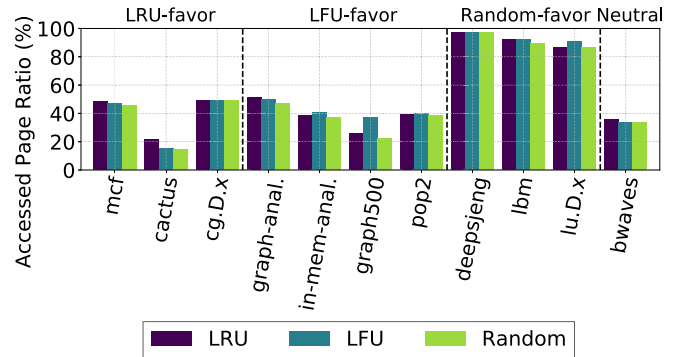


Fig. 8. Average accessed page ratio of workloads.

memory hit ratio shows the strongest correlation (-0.8194) and the smallest p-value (1.0505×10^{-11}). The accessed page ratio has a p-value lower than 0.01, showing that the feature has a statistically significant correlation. However, the absolute value of the correlation coefficient is smaller than the fast memory hit ratio's. On the other hand, the page migration stability does not have a statistically meaningful relationship with the performance (p-value = 0.8355).

A Feature for Random-Favor Workloads. We find that having a high accessed page ratio is a hint for workloads to favor the random migration policy. As we have presented in Section 3.2, random-favor workloads access memory with a low locality, and they have a huge memory footprint that exceeds the fast memory size. As a result, random-favor workloads have a higher average accessed page ratio compared to the other workloads. Fig. 8 presents the average accessed page ratio of workloads. Random-favor workloads have average accessed page ratios higher than 80 percent.

4.2 Adaptive Page Migration Policy Selection

AMP adaptively selects a page migration policy between LRU, LFU, and Random using the features that we have analyzed in the previous subsection. Algorithm 1 describes the page migration policy decision of AMP. AMP classifies a workload as random-favor if the accessed page ratio of the workload exceeds the fast memory ratio significantly. The fast memory ratio is defined as the number of pages in fast memory divided by the number of total pages. The insight behind this heuristic is that a workload with a huge working set that exceeds the cache size may experience a thrashing. If the accessed page ratio exceeds the fast memory ratio by 20 percent, AMP chooses the random migration policy. The threshold is empirically set.

Otherwise, AMP selects a page migration policy between LRU and LFU. According to the feature analysis, the fast memory hit ratio and accessed page ratio have a strong relationship with the performance. Between the two features, we choose the fast memory hit ratio because its correlation is stronger than the other. AMP tracks the fast memory hit ratios of LRU and LFU simultaneously and chooses a policy that has a higher average fast memory hit ratio. The key challenge in tracking the fast memory hit ratios is that only one page migration policy can be applied to physical tiered memory. We overcome this problem by emulating a page migration policy. Fig. 9 illustrates how AMP obtains the fast memory hit ratios of both policies simultaneously. AMP applies the page migration policy with the higher average

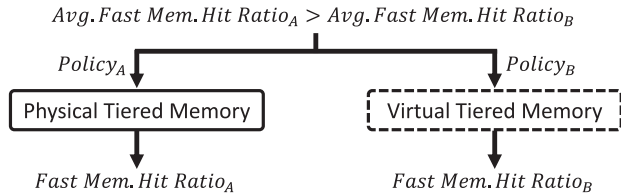


Fig. 9. Adaptive page migration policy selection for non-random-favor workloads.

fast memory hit ratio to physical tiered memory. At the same time, pages are migrated *virtually* with the other policy. For example, if LRU has a higher average fast memory hit ratio, LRU is applied to physical pages. At the same time, LFU is emulated *virtually* in the kernel without page migrations. As a result, the kernel can collect the fast memory hit ratios of both policies. If LFU turns out to have a higher average fast memory hit ratio, AMP applies LFU to the physical tiered memory in the next turn. AMP uses the moving average of fast memory hit ratios, and the window size is 36-epoch. We empirically find that the moving average can select a page migration policy accurately and stably.

Algorithm 1. Adaptive Page Migration Policy Selection

```

1: if accessed_page_ratio > (fast_memory_ratio + 20%) then
2:   repl_policy = Random
3: else
4:   if LRU_hit_ratio_avg > LFU_hit_ratio_avg then
5:     repl_policy = LRU
6:   else
7:     repl_policy = LFU

```

Cost of Switching Policies. Switching the policy from one to another has a negligible cost, as it affects only which pages need to be in the fast memory. To choose the right policy, the proposed component tracks the hotness, fast memory hit ratio, and other information. However, this tracking computation is done in the background during each interval. An indirect cost is that some pages in the fast memory which were promoted by the old policy, may no longer be hot ones with the new policy. Therefore, gradually the pages are evicted, and new pages are migrated by the new policy.

Although the policy switching itself does not have any significant direct overheads, the extra information must be tracked and maintained to choose the right policy. It has some memory capacity overheads and computational costs. The spatial cost includes the per-page metadata to track hotness (4B age, 8B access history, 4B access frequency), per-page virtual page location (4B) to simulate page migrations, and per-policy features (fast memory hit ratio and accessed page ratio, 4B each). The majority of computational costs to simulate the other policy in the background without actual page migrations. This cost occurs when it sorts pages to find the relative hotness of pages. Note that the computation occurs in the background during each time interval, not during the policy switching.

5 IMPLEMENTATION

We implement AMP in a Linux system. The implementation of AMP spans from the kernel to the user-level. The kernel tracks the age and access frequency of pages and provides

several options for page migration policies. Additionally, features such as the fast memory hit ratio and accessed page ratio are collected in the kernel. AMP is built on the memory control groups in the kernel (memcgs). We assume that the processes in a memcg have the same preferences on page migration policies. Therefore, memcg is the basic unit of page migration policy decisions. The fast memory ratio can be set for each memcg.

The user-level controller periodically requests the scan of pages. The request is sent to the kernel by writing to a file under `sysfs`. On the request, the kernel scans all pages in the memcg and checks the accessed bits of the pages. Page age and access frequency are updated using the accessed bits. The accessed bits of pages are checked using the `page_is_idle` function [47]. After that, the accessed bit of the page is unset using the `set_page_idle` function to check further accesses to the page.

The user-level controller migrates pages between fast memory and slow memory by requesting to the kernel. The user-level controller sets the page migration policy, and the kernel applies the page migration policy. AMP chooses the random replacement policy if the accessed page ratio exceeds a predefined threshold. Otherwise, the policy with a higher average fast memory hit ratio is selected between LRU and LFU. The kernel reports the fast memory hit ratio of both policies to the user-level controller. The user-level controller collects the fast memory hit ratios and calculates the moving averages of fast memory hit ratios.

Page migration requires exchanges of pages between two nodes. In the native Linux, page exchanges involve redundant page (de)allocations, which causes the performance overhead. The recently proposed optimization can eliminate the overhead [13]. The proposed optimization exchanges two pages by changing the mappings and exchanging the contents of pages without (de)allocating pages. We apply the kernel patch [46] to reduce the performance overhead of page migrations.

6 EVALUATION

6.1 Experiment Setup

We evaluate AMP on a Linux system. The system runs as a two-socket QEMU virtual machine to emulate a tiered memory system. The system is composed of fast and slow memory nodes. The fast memory node has CPU cores, and its memory is allocated from the normal DRAM. The slow memory node does not have CPU cores, and its memory is allocated from a bandwidth-throttled DRAM. We throttle the memory bandwidth using power throttling [48], and we saturate the memory bandwidth using `membw` [49] to meet the reported latency (346ns) of Optane DC [2]. Table 2 describes the evaluation system configurations.

6.2 Performance of AMP

Fig. 10 shows the normalized performance of workloads with various page migration policies. We compare the performance of workloads with the modified LRU lists, LRU, LFU, Random, and AMP. For the modified LRU lists, we run the experiments with 4 KB base pages and 2 MB transparent huge pages, respectively. The suffix in the legend shows the page size. `b` stands for base pages, and `h` means huge pages. For the other configurations, 2 MB transparent huge pages are used.

TABLE 2
System Configurations

Intel Xeon Dual Socket System	
OS & Kernel	Ubuntu 18.04.2 - kernel 4.15.0
Processors	2-socket E5-2630 v4
Memory	DDR4 - 2133MHz
Fast Memory Latency	78ns
Fast Memory BW	32 GB/s
Slow Memory Latency (Emulated)	359ns
Slow Memory BW (Emulated)	5.8 GB/s

TABLE 3
Huge Page Allocation Ratio in Anonymous Pages

Workload Name	2MB Ratio	Workload Name	2MB Ratio
mcf	94%	pop2	77%
cactus	92%	deepsjeng	97%
cg.D.x	100%	lbn	94%
graph-analytics	97%	lu.D.x	99%
in-mem-analytics	96%	bwaves	98%
graph500	98%	Geomean	95%

Table 3 shows the portion of huge page allocation ratio of each workload, presenting 95 percent on average. The performance is the reverse of the execution time, and the performance is normalized to the 100 percent fast memory ratio. In this experiment, we set the fast memory ratio to 50 percent. On average, AMP can achieve 10.9, 6.4, 17.6 percent higher performance compared to LRU, LFU, and Random, respectively.

Fig. 11a shows the temporal change of average fast memory hit ratios, and Fig. 11b presents the timeline of page migration policy selections. *cactus*, *graph500*, and *lbn* favor LRU, LFU, and Random, respectively. Overall, the preferred page migration policy has a higher average fast memory hit ratio during the execution time except for *lbn*. For *cactus*, LRU shows the higher average fast memory

hit ratio. Therefore, *cactus* selects LRU except for the warming-up stage. For *graph500*, LRU has a higher average fast memory hit ratio in the initial stage. Therefore, *graph500* chooses LRU at first. As time goes by, the average fast memory hit ratio of LFU beats the LRU's. After the point, *graph500* chooses LFU. On the other hand, *lbn* does not show any difference between the fast memory hit ratios of LRU and LFU because it favors Random. Fig. 11 shows that AMP can choose a page migration policy using the average fast memory hit ratios.

6.3 Page Migration Policy Selection Ratio

We measure the AMP's page migration policy selection ratio for each workload. Fig. 12 presents the page migration

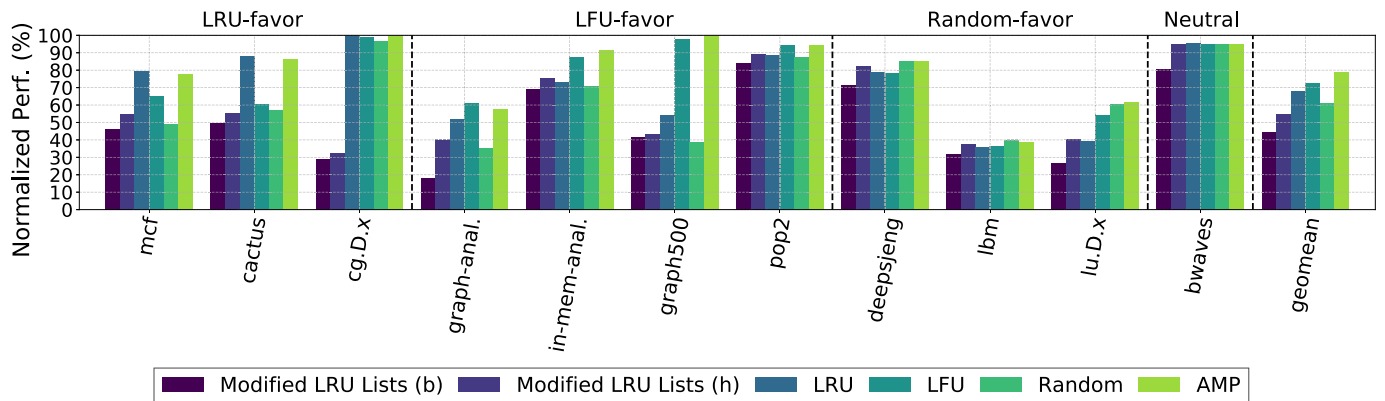


Fig. 10. Normalized performance of workloads with various page migration policies.

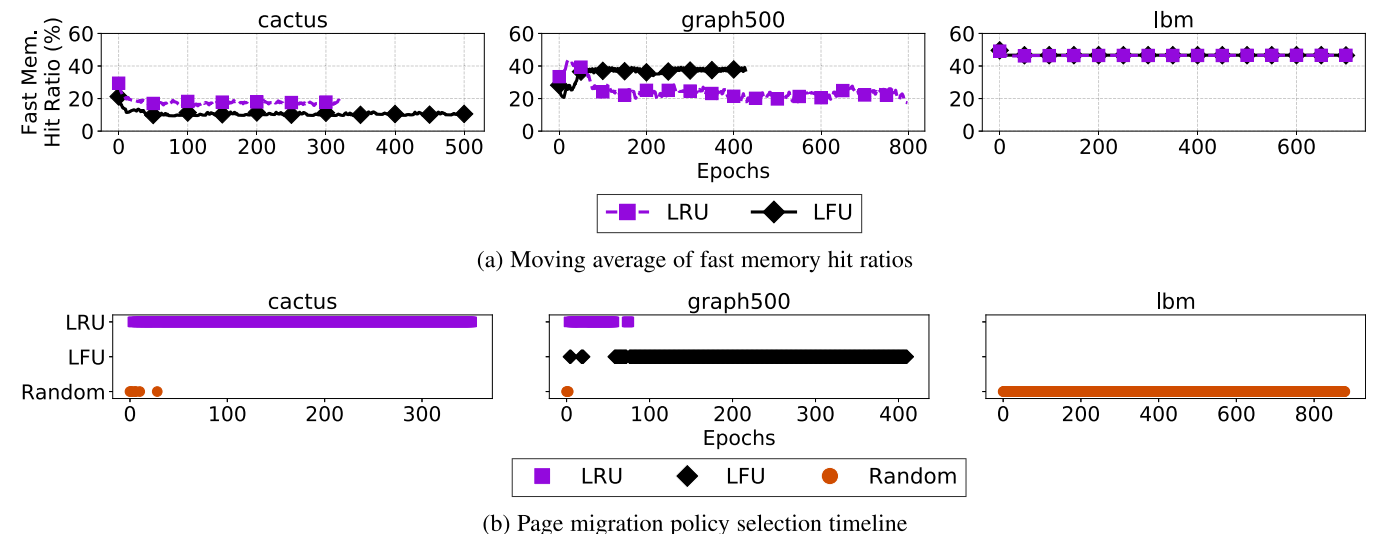


Fig. 11. Temporal change of fast memory hit ratios and page migration policy selection timeline.

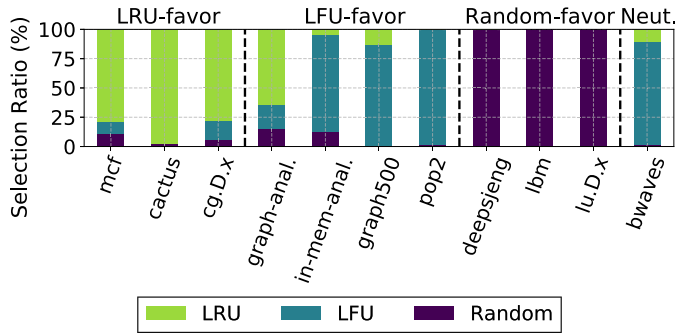


Fig. 12. Page migration policy selection ratio.

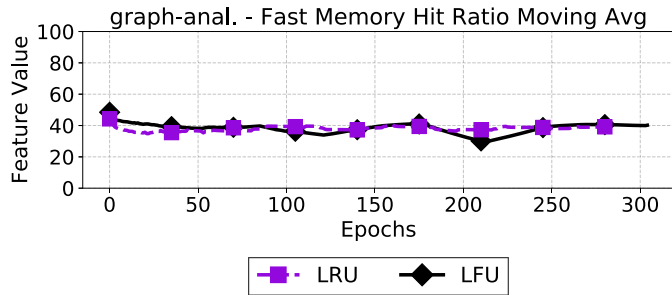


Fig. 13. Moving averages of fast memory hit ratios (graph-analytics).

policy selection ratio for each workload. Overall, AMP can choose a preferred page migration policy. For most workloads, the preferred page migration policy shows a selection ratio higher than 80 percent. However, *graph-analytics* presents the high selection ratio of LRU, although it prefers LFU when a policy is chosen statically. It is because the average fast memory hit ratios of LRU and LFU do not show a meaningful gap during most of the execution time. Fig. 13 shows the temporal change of the average fast memory hit ratios of *graph-analytics*, and the figure presents the negligible gap between LRU and LFU. Therefore, choosing LRU is okay for *graph-analytics*. Another reason for a workload to choose a non-favored page migration policy is the change in page migration policy preference during the execution, as we have shown in *graph500* in the previous subsection.

6.4 Performance of AMP in Consolidated Environments

One of the advantages of AMP is that it can apply a different page migration policy for each memcg, simultaneously.

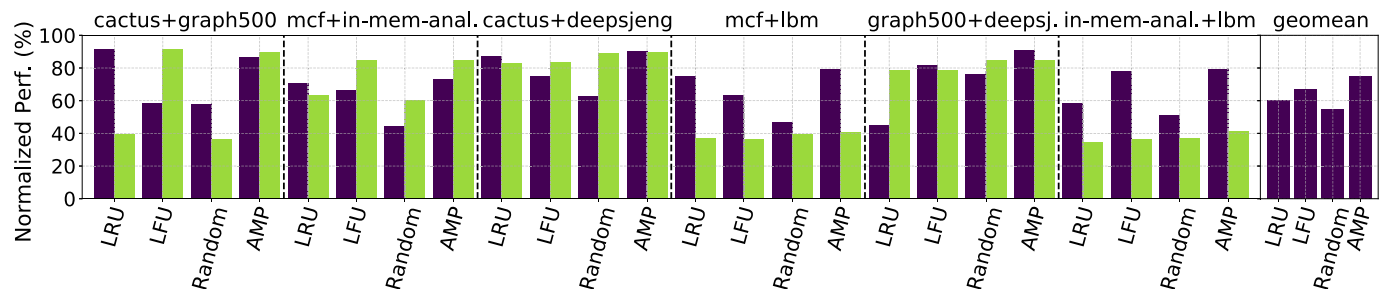


Fig. 14. Normalized performance of workload mixes with various page migration policies (consolidated).

TABLE 4
Page Migration Policy Preferences Summary

LRU-favor	mcf, cactus, cg.D.x
LFU-favor	graph-anal., in-mem-anal., graph500, pop2
Random-favor	deepsjeng, lbm, lu.D.x
Neutral	bwaves

Even though a set of workloads have different favors on page migration policies, AMP can offer the preferred page migration policy for each workload. In this subsection, we evaluate the performance of AMP when multiple workloads are running simultaneously in different memcgs. We run six workload mixes that are combinations of workloads with different page migration policy preferences. The fast memory ratio is set to 50 percent of the working sets.

Fig. 14 shows the normalized performance of workload mixes. In the figure, there are seven groups of bar graphs. The first six groups represent workload mixes, and the corresponding workload mix is shown above the figure. The workloads' preferences on page migration policies are summarized in the Table 4. The first six groups are composed of grouped bar graphs. The first bar in a group shows the performance of the first workload within a mix, and the second bar presents the performance of the second workload in a mix. We use the same definition for the performance that we use in Section 6.2. The last group of bars shows the geomean performance of page migration policies. Each workload shows the best performance with the preferred page migration policy. Static policies such as LRU, LFU, and Random cannot offer the best performance for both workloads in a mix at the same time. AMP can offer the preferred policy for both workloads in a mix, achieving 14.7, 8.2, 20.4 percent higher geomean performance than LRU, LFU, and Random, respectively.

6.5 Sensitivity to the Slow Memory Ratio

We evaluate the sensitivity of AMP to the fast memory ratio. In the previous experiments, we set the fast memory ratio to 50 percent of workloads' working set. In this section, we show that AMP performs better than the other page migration policies with various fast memory ratios. We measure the performance of workloads with 30, 50, and 70 percent fast memory ratios. We use the same workloads, page migration policies, and the same definition for the performance that we use in Section 6.2. Fig. 16 presents the normalized performance of workloads with various page migration policies. We show the geomean performance of workloads. Overall, AMP performs better than the other

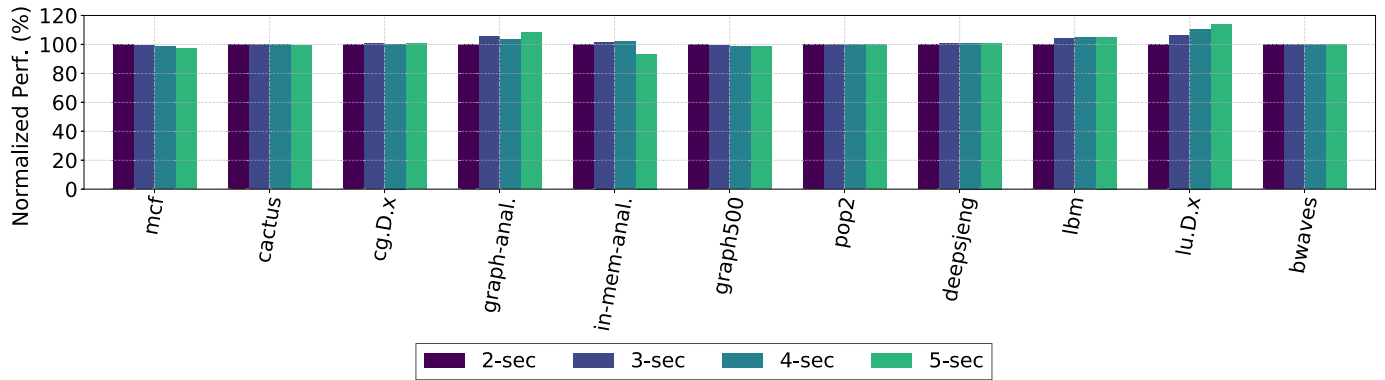


Fig. 15. Normalized performance of workloads with various page migration intervals.

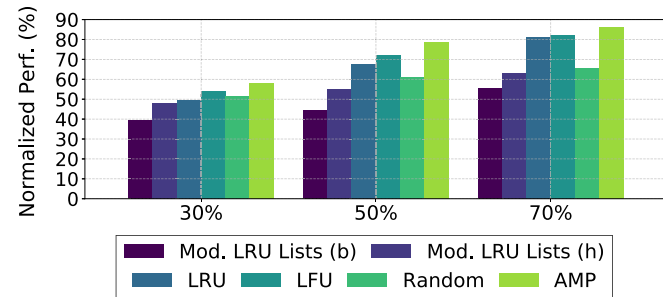


Fig. 16. Sensitivity to the fast memory ratio.

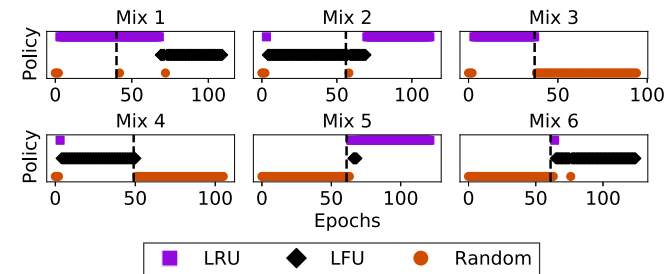


Fig. 17. Timeline of page migration policy selections.

page migration policies regardless of the fast memory ratio. AMP shows 4.0, 6.4, 4.1 percent higher performance than LFU at 30, 50, 70 percent fast memory ratio, respectively.

6.6 Responsiveness to the Phase Transitions

AMP works well for most macro-benchmarks with a gradual transition between phases. A phase is an interval of execution where the preference on page migration policies is the same. However, if a phase changes abruptly, AMP may experience a delay until updating the page migration policy selection because it uses the moving average of fast memory hit ratios. In this subsection, we measure the delay until AMP learns the changes in the page migration policy preference on abrupt phase transitions.

We use three types of synthetic benchmarks.

LRU-favor has a strided memory access pattern with four same-sized working sets. Each working set is sequentially accessed for four seconds.

LFU-favor has two same-sized working sets. One is a frequently-accessed hot working set, and the other is an infrequently-accessed cold working set. The infrequently-

TABLE 5
Synthetic Benchmark Mixes

Mix 1	LRU-favor → LFU-favor
Mix 2	LFU-favor → LRU-favor
Mix 3	LRU-favor → Random-favor
Mix 4	LFU-favor → Random-favor
Mix 5	Random-favor → LRU-favor
Mix 6	Random-favor → LFU-favor

TABLE 6
The Number of Epochs Until AMP Learns the Changes in the Preference on Page Migration Policy

Mix 1	27-epoch	Mix 2	17-epoch	Mix 3	2-epoch
Mix 4	2-epoch	Mix 5	3-epoch	Mix 6	4-epoch

accessed cold working set is divided into 16 subsets, showing a strided memory access pattern with a low access frequency.

Random-favor allocates memory, and it randomly accesses all pages. We compose six types of workload mixes with synthetic benchmarks. Table 5 presents the execution order of workload mixes.

Fig. 17 illustrates the timeline of page migration policy selections for each workload mix. Each workload starts with Random in the warming-up stage, and it selects the preferred page migration policy after the warming up. The dashed vertical line shows the transition point, where the first workload exits and the second workload starts execution. Except for Mix 1, AMP can follow the preferred page migration policy within a relatively short delay. Table 6 summarizes the delay until AMP finds the favored page migration policy after a phase transition. AMP shows the low responsiveness in Mix 1 because of the stale fast memory hit ratios kept in moving average. LRU shows the high fast memory hit ratio in the first workload of Mix 1, and it takes time for AMP to learn the changes in policy preferences. Please note that this kind of abrupt phase transitions are not common in the real world.

6.7 Sensitivity to the Page Migration Interval

Page migration interval is one of the important parameters in migrating pages in tiered memory systems. There is a trade-off between responsive migration of hot pages to fast

memory and page migration cost. In the previous subsections, the page migration interval is set to five seconds empirically. In this subsection, we evaluate the sensitivity of AMP's performance to the page migration interval with shorter intervals than five seconds. Fig. 15 presents the performance of workloads with AMP with various page migration intervals. The performance is normalized to the performance with the shortest page migration interval, 2-second. Most workloads show similar performance regardless of the page migration interval. *1u.D.x* presents 13.8 percent higher performance with 5-second page migration interval compared to 2-second interval's. This is because the memory access pattern of *1u.D.x* is random, and most pages are actively accessed. Therefore, frequent page migrations add mere performance overhead without increasing the fast memory hit ratio. *In-memory-analytics* presents the higher performance with the intervals shorter than 5-second, implying that the shorter page migration interval can capture the dynamic nature of hot working set of *in-memory-analytics*'s. To summarize, there are some workloads that are sensitive to page migration interval. However, there is a little gain in shortening page migration interval, on average.

7 DISCUSSION

Low-Overhead Hotness Tracking Mechanisms. Having a low-cost hotness tracking mechanism is important in identifying the hotness of pages. Checking the accessed bits of pages at a low frequency is one of the solutions to achieve the low overhead [19]. Although this is a viable solution to identify swap page candidates, it cannot be applied to migrating pages in tiered memory because of its low resolution. Alternatively, dynamically adjusting the unit of hotness tracking is a solution to reduce the performance overhead [50]. By identifying groups of pages that have similar page hotness, the number of accessed bit checks can be reduced. AMP achieves the low overhead in hotness tracking by checking the accessed bits at the granularity of huge pages.

Comparison With HW-Based Migration Mechanisms. The target memory system of this study is a tiered memory system where software is responsible for managing page locations between memory tiers. Hardware data migration mechanisms [51], [52], [53] have an advantage in offering software-transparent fine-grained data migration. However, the hardware mechanisms require modifications to memory controllers, which needs support from hardware vendors. The proposed software mechanism and policies can be applied without any support from hardware vendors, which goes well with the current data centers.

The hardware techniques [51], [52], [53] are intended for a hybrid memory system with a relatively small 3D stacked DRAM (fast memory) backed by the conventional DRAM (slow memory). Therefore, the fast memory capacity is smaller than what is used in general tiered memory where DRAM and NVM are combined. Thanks to the small fast memory capacity of 3D stacked DRAM, the HW approach maintains an extra layer of mapping between the fast and slow memory spaces, instead of using page tables. With the HW maintained mapping table, it is possible to migrate data at smaller granularity than pages in a nimble way.

Those HW approaches access the mapping table slightly differently. Some approaches cache the HW-maintained mapping table in on-chip SRAM for fast access [51], while the other approaches look up the fast memory for mapping information whenever an LLC miss occurs [53]. Both of the mechanisms were possible since they are designed for the 3D stacked DRAM, which was assumed to be faster than DRAM for fast lookups of the mapping table, and to have a small capacity so that the mapping table can be efficiently cached in SRAM. However, the fine-grained HW mapping tables may not be scalable enough to cover the combined capacity of DRAM and NVM, and accessing DRAM first for mapping information for every LLC miss will slow down memory access times significantly. In addition, the HW mapping table has a limited associativity to reduce the size as much as possible. Therefore, the memory management is severely restricted, and there are only a few locations data can be stored either fast or slow memory.

Cost of TLB Shootdowns. When a page table entry is updated for a process, the Linux kernel sends Inter-Process Interrupts (IPIs) to the cores running threads of the same process for TLB shootdowns. Once an IPI arrives, the receiving core invokes the kernel to execute TLB invalidation. IPIs are sent to the cores running threads with the same address space, regardless of whether pages are actually shared or not [54], [55]. Note that the OS kernel does not track which pages are shared by what threads within a process. Therefore, even with a 4 KB base page, an update of a PTE will send IPIs to the cores running the other threads of the same address space, even if the other threads do not access the page and the cores do not have the affected PTE in their TLBs. The majority of the shutdown cost is for initiating and responding to IPIs, regardless of TLB hits or misses during the TLB invalidation. Therefore, huge pages do not increase the occurrences of TLB shootdowns by data sharing, compared to base pages.

Instead, using huge pages can potentially reduce TLB shutdown occurrences. If the entire region of a huge page is accessed, a single shutdown can migrate a 2 MB page. With 4 KB base pages, 512 shutdowns can occur in the worst case for a migration of the same 2 MB region. Note that shutdown IPIs will be sent to the cores running all threads in the same process, even if a PTE of any base or huge page is updated. In addition, using huge pages significantly reduces TLB misses for many workloads.

Partitioning Fast Memory. Although this study assumes a case where each workload uses fast memory exclusively, sharing fast memory can be preferred in a batch execution environment. Sharing fast memory between multiple workloads is similar to partitioning CPU caches between multiple processes. This problem has been studied for several decades in the computer architecture community. The insights from cache partitioning studies can be applied to sharing fast memory in tiered memory systems. Cache partitioning has been studied to allocate caches between multiple processes to minimize the miss rate and maximize throughput [56], [57], to guarantee the fairness between applications [56], [58], [59], [60], [61], [62], [63], and to protect the latency-sensitive jobs from batch jobs [64], [65], [66], [67]. Constructing miss rate curves or utility curves can give users a hint to allocate fast memory between multiple

workloads [57], [68], [69], [70], [71], [72]. If users already know the utility curves of workloads, an auction can be used to allocate fast memory between workloads [73], [74].

8 CONCLUSION

Memory systems are adopting memories with different latency and bandwidth, comprising a tiered memory system. Page migration policies migrate pages to utilize fast memory with hot pages. We find that workloads have diverse preferences on page migration policies. We analyze the reason behind the various preferences on policies, and we find the relationship between features and performance of workloads. Based on the analysis, we propose AMP, which adaptively selects a page migration policy between LRU, LFU, and Random. AMP can estimate the fast memory hit ratio of a page migration policy by emulating the policy without page migrations. AMP can achieve 10.9, 6.4, 17.6percent higher performance than LRU, LFU, and Random, respectively. The source code is available at <https://github.com/casys-kaist/AMP>.

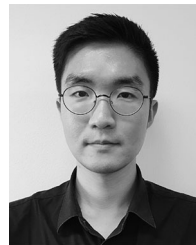
ACKNOWLEDGMENTS

The authors would like to thank the anonymous reviewers for their insightful feedbacks and comments. This work was supported in part by the National Research Foundation of Korea under Grant NRF-2019R1A2B5B01069816 and in part by Institute for Information & communications Technology Promotion under Grant IITP-2017-0-00466.

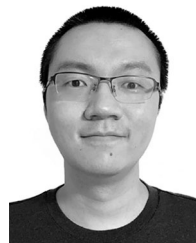
REFERENCES

- [1] Intel® Optane™ DC Persistent Memory. Accessed: Jan. 13, 2020. [Online]. Available: <https://www.intel.com/content/www/us/en/architecture-and-technology/optane-dc-persistent-memory.html>
- [2] J. Izraelevitz *et al.*, “Basic performance measurements of the Intel Optane DC persistent memory module,” *CoRR*, vol. abs/1903.05714, 2019. [Online]. Available: <http://arxiv.org/abs/1903.05714>
- [3] R. L. Mattson, J. Gecsei, D. R. Slutz, and I. L. Traiger, “Evaluation techniques for storage hierarchies,” *IBM Syst. J.*, vol. 9, no. 2, pp. 78–117, 1970.
- [4] D. Shasha and T. Johnson, “2Q: A low overhead high performance buffer management replacement algorithm,” in *Proc. 20th Int. Conf. Very Large Databases*, 1994, pp. 439–450.
- [5] D. Lee *et al.*, “LRFU: A spectrum of policies that subsumes the least recently used and least frequently used policies,” *IEEE Trans. Comput.*, vol. 50, no. 12, pp. 1352–1361, Dec. 2001.
- [6] I. Ari, A. Amer, R. B. Gramacy, E. L. Miller, S. A. Brandt, and D. D. Long, “ACME: Adaptive caching using multiple experts,” in *Proc. Workshop Distrib. Data Struct.*, 2002, pp. 143–158.
- [7] S. Jiang and X. Zhang, “LIRS: An efficient low inter-reference recency set replacement policy to improve buffer cache performance,” in *Proc. ACM SIGMETRICS Int. Conf. Meas. Model. Comput. Syst.*, 2002, pp. 31–42.
- [8] N. Megiddo and D. S. Modha, “ARC: A self-tuning, low overhead replacement cache,” in *Proc. 2nd USENIX Conf. File Storage Technol.*, 2003, pp. 115–130.
- [9] S. Bansal and D. S. Modha, “CAR: Clock with adaptive replacement,” in *Proc. 3rd USENIX Conf. File Storage Technol.*, 2004, pp. 187–200.
- [10] S. Jiang, F. Chen, and X. Zhang, “CLOCK-Pro: An effective improvement of the CLOCK replacement,” in *Proc. USENIX Annu. Tech. Conf.*, 2005, pp. 323–336.
- [11] M. K. Qureshi, A. Jaleel, Y. N. Patt, S. C. Steely, and J. Emer, “Adaptive insertion policies for high performance caching,” in *Proc. 34th Int. Symp. Comput. Architect.*, 2007, pp. 381–391.
- [12] N. Agarwal and T. F. Wenisch, “Thermostat: Application-transparent page management for two-tiered main memory,” in *Proc. 22nd Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2017, pp. 631–644.
- [13] Z. Yan, D. Lustig, D. Nellans, and A. Bhattacharjee, “Nimble page management for tiered memory systems,” in *Proc. 24th Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2019, pp. 331–345.
- [14] Intel memory drive technology. Accessed: Jan. 13, 2020. [Online]. Available: <https://www.intel.com/content/www/us/en/software/intel-memory-drive-technology.html>
- [15] J. Gu, Y. Lee, Y. Zhang, M. Chowdhury, and K. G. Shin, “Efficient memory disaggregation with INFINISWAP,” in *Proc. 14th USENIX Symp. Netw. Syst. Des. Implementation*, 2017, pp. 649–667.
- [16] M. K. Aguilera *et al.*, “Remote regions: A simple abstraction for remote memory,” in *Proc. USENIX Annu. Tech. Conf.*, 2018, pp. 775–787.
- [17] V. Nitu, B. Teabe, A. Tchana, C. Isci, and D. Hagimont, “Welcome to Zombieland: Practical and energy-efficient memory disaggregation in a datacenter,” in *Proc. 13th Eur. Conf. Comput. Syst.*, 2018, pp. 1–12.
- [18] K. Koh, K. Kim, S. Jeon, and J. Huh, “Disaggregated cloud memory with elastic block management,” *IEEE Trans. Comput.*, vol. 68, no. 1, pp. 39–52, Jan. 2019.
- [19] A. Lagar-Cavilla *et al.*, “Software-defined far memory in warehouse-scale computers,” in *Proc. 24th Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2019, pp. 317–330.
- [20] K. Keeton, “The machine: An architecture for memory-centric computing,” in *Proc. Int. Workshop Runtime Operating Syst. Supercomput.*, 2015, Art. no. 1.
- [21] Transparent hugepages. Accessed: Jan. 13, 2020. [Online]. Available: <https://lwn.net/Articles/359158/>
- [22] J. Navarro, S. Iyer, P. Druschel, and A. Cox, “Practical, transparent operating system support for superpages,” in *Proc. 5th Symp. Operating Syst. Des. Implementation*, 2002, pp. 89–104.
- [23] Y. Kwon, H. Yu, S. Peter, C. J. Rossbach, and E. Witchel, “Coordinated and efficient huge page management with Ingens,” in *Proc. 12th Symp. Operating Syst. Des. Implementation*, 2016, pp. 705–721.
- [24] A. Panwar, A. Prasad, and K. Gopinath, “Making huge pages actually useful,” in *Proc. 23rd Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2018, pp. 679–692.
- [25] A. Panwar, S. Bansal, and K. Gopinath, “HawkEye: Efficient fine-grained OS support for huge pages,” in *Proc. 24th Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2019, pp. 347–360.
- [26] F. J. Corbato, “A paging experiment with the multics system,” MIT Press, Cambridge, MA, USA, MIT Project MAC Rep. MAC-M-384, May 1968.
- [27] M. Kampe, P. Stenstrom, and M. Dubois, “Self-correcting LRU replacement policies,” in *Proc. 1st Conf. Comput. Front.*, 2004, pp. 181–191.
- [28] E. Teran, Y. Tian, Z. Wang, and D. A. Jiménez, “Minimal disturbance placement and promotion,” in *Proc. 22nd Int. Symp. High Perform. Comput. Architect.*, 2016, pp. 201–211.
- [29] V. Gupta, M. Lee, and K. Schwan, “HeteroVisor: Exploiting resource heterogeneity to enhance the elasticity of cloud platforms,” in *Proc. 11th Int. Conf. Virt. Execution Environ.*, 2015, pp. 79–92.
- [30] G. Glass and P. Cao, “Adaptive page replacement based on memory reference behavior,” in *Proc. ACM SIGMETRICS Int. Conf. Meas. Model. Comput. Syst.*, 1997, pp. 115–126.
- [31] Y. Smaragdakis, S. Kaplan, and P. Wilson, “EELRU: Simple and effective adaptive page replacement,” in *Proc. ACM SIGMETRICS Int. Conf. Meas. Model. Comput. Syst.*, 1999, pp. 122–133.
- [32] J. M. Kim *et al.*, “A low-overhead high-performance unified buffer management scheme that exploits sequential and looping references,” in *Proc. 4th Symp. Operating Syst. Des. Implementation*, 2000, Art. no. 9.
- [33] J. T. Robinson and M. V. Devarakonda, “Data cache management using frequency-based replacement,” in *Proc. ACM SIGMETRICS Int. Conf. Meas. Model. Comput. Syst.*, 1990, pp. 134–142.
- [34] E. J. O’neil, P. E. O’neil, and G. Weikum, “The LRU-K page replacement algorithm for database disk buffering,” in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 1993, pp. 297–306.
- [35] Y. Zhou, J. Philbin, and K. Li, “The multi-queue replacement algorithm for second level buffer caches,” in *Proc. USENIX Annu. Tech. Conf.*, 2001, pp. 91–104.
- [36] D. Lee *et al.*, “On the existence of a spectrum of policies that subsumes the least recently used (LRU) and least frequently used (LFU) policies,” in *Proc. ACM SIGMETRICS Int. Conf. Meas. Model. Comput. Syst.*, 1999, pp. 134–143.
- [37] G. Vietri *et al.*, “Driving cache replacement with ML-based LeCar,” in *Proc. USENIX Workshop Hot Topics Storage File Syst.*, 2018, Art. no. 3.

- [38] SPEC CPU 2017. Accessed: Jan. 13, 2020. [Online]. Available: <https://www.spec.org/cpu2017/>
- [39] M. Ferdman *et al.*, "Clearing the clouds: A study of emerging scale-out workloads on modern hardware," in *Proc. 17th Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2012, pp. 37–48.
- [40] NAS parallel benchmarks. Accessed: Jan. 13, 2020. [Online]. Available: <https://www.nas.nasa.gov/publications/npb.html>
- [41] Graph500. Accessed: Jan. 13, 2020. [Online]. Available: <https://graph500.org/>
- [42] C.-K. Luk, R. Muth, H. Patil, R. Weiss, P. G. Lowney, and R. Cohn, "Profile-guided post-link stride prefetching," in *Proc. 16th Int. Conf. Supercomput.*, 2002, pp. 167–178.
- [43] H. Al-Sukhni *et al.*, "The design of cost-effective stride-prefetching for modern processors," Dept. Elect. Comput. Eng., Tech. Rep., Freescale Semiconductor, Inc. Austin, TX, Univ. Colorado Boulder,
- [44] R. Panda and L. K. John, "HALO: A hierarchical memory access locality modeling technique for memory system explorations," in *Proc. 32nd Int. Conf. Supercomput.*, 2018, pp. 118–128.
- [45] C. Takahashi, *et al.*, "Empirical study for optimization of power-performance with on-chip memory," in *Proc. Int. Symp. High-Perform. Comput.*, 2005, pp. 466–479.
- [46] Nimble page management for tiered memory systems - GitHub repository. Accessed: Jan. 13, 2020. [Online]. Available: https://github.com/ysarch-lab/nimble_page_management_aspos_2019/blob/5d503c456f1eeced24e4723ef758ce2c38db1ae0/mm/memory_manage.c#L777
- [47] Idle page tracking. Accessed: Jan. 13, 2020. [Online]. Available: https://www.kernel.org/doc/html/latest/admin-guide/mm/idle_page_tracking.html
- [48] Intel® Xeon Processor E5 v4 Product Family. Accessed: Jan. 13, 2020. [Online]. Available: <https://www.intel.com/content/dam/www/public/us/en/documents/datasheets/xeon-e5-v4-datasheet-vol-2.pdf>
- [49] Intel-cmt-cat. Accessed: Jan. 13, 2020. [Online]. Available: <https://github.com/intel/intel-cmt-cat>
- [50] S. Park, Y. Lee, and H. Y. Yeom, "Profiling dynamic data access patterns with controlled overhead and quality," in *Proc. 20th Int. Middleware Conf. Ind. Track*, 2019, pp. 1–7.
- [51] J. Sim, A. R. Alameldeen, Z. Chishti, C. Wilkerson, and H. Kim, "Transparent hardware management of stacked DRAM as part of memory," in *Proc. 47th Int. Symp. Microarchitect.*, 2014, pp. 13–24.
- [52] J. B. Kotra, H. Zhang, A. R. Alameldeen, C. Wilkerson, and M. T. Kandemir, "CHAMELEON: A dynamically reconfigurable heterogeneous memory system," in *Proc. 51st Int. Symp. Microarchitect.*, 2018, pp. 533–545.
- [53] C. C. Chou, A. Jaleel, and M. K. Qureshi, "CAMEO: A two-level memory organization with capacity of main memory and flexibility of hardware-managed cache," in *Proc. 47th Int. Symp. Microarchitect.*, 2014, pp. 1–12.
- [54] N. Amit, A. Tai, and M. Wei, "Don't shoot down TLB shootdowns!" in *Proc. 15th Eur. Conf. Comput. Syst.*, 2020, pp. 1–14.
- [55] M. K. Kumar *et al.*, "LATR: Lazy translation coherence," in *Proc. 23rd Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2018, pp. 651–664.
- [56] L. R. Hsu, S. K. Reinhardt, R. Iyer, and S. Makineni, "Communist, utilitarian, and capitalist cache policies on CMPs: Caches as a shared resource," in *Proc. 15th Int. Conf. Parallel Architect. Compilation Techn.*, 2006, pp. 13–22.
- [57] M. K. Qureshi and Y. N. Patt, "Utility-based cache partitioning: A low-overhead, high-performance, runtime mechanism to partition shared caches," in *Proc. 39th Int. Symp. Microarchitect.*, 2006, pp. 423–432.
- [58] S. Kim, D. Chandra, and Y. Solihin, "Fair cache sharing and partitioning in a chip multiprocessor architecture," in *Proc. 13th Int. Conf. Parallel Architect. Compilation Techn.*, 2004, pp. 111–122.
- [59] T. Y. Yeh and G. Reinman, "Fast and fair: Data-stream quality of service," in *Proc. Int. Conf. Compilers Architect. Synthesis Embedded Syst.*, 2005, pp. 237–248.
- [60] J. Chang and G. S. Sohi, "Cooperative cache partitioning for chip multiprocessors," in *Proc. 25th Int. Conf. Supercomput.*, 2007, pp. 402–412.
- [61] X. Wang and J. F. Martínez, "ReBudget: Trading off efficiency vs. fairness in market-based multicore resource allocation via runtime budget reassignment," in *Proc. 21st Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2016, pp. 19–32.
- [62] V. Selfa, J. Sahuquillo, L. Eeckhout, S. Petit, and M. E. Gómez, "Application clustering policies to address system fairness with Intel's cache allocation technology," in *Proc. 26th Int. Conf. Parallel Architect. Compilation Techn.*, 2017, pp. 194–205.
- [63] J. Park, S. Park, and W. Baek, "CoPart: Coordinated partitioning of last-level cache and memory bandwidth for fairness-aware workload consolidation on commodity servers," in *Proc. 14th Eur. Conf. Comput. Syst.*, 2019, pp. 1–14.
- [64] H. Cook, M. Moreto, S. Bird, K. Dao, D. A. Patterson, and K. Asanovic, "A hardware evaluation of cache partitioning to improve utilization and energy-efficiency while preserving responsiveness," in *Proc. 40th Int. Symp. Comput. Architect.*, 2013, pp. 308–319.
- [65] H. Kasture and D. Sanchez, "Ubik: Efficient cache sharing with strict QoS for latency-critical workloads," in *Proc. 19th Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2014, pp. 729–742.
- [66] D. Lo, L. Cheng, R. Govindaraju, P. Ranganathan, and C. Kozyrakis, "Heracles: Improving resource efficiency at scale," in *Proc. 42nd Int. Symp. Comput. Architect.*, 2015, pp. 450–462.
- [67] H. Zhu and M. Erez, "Dirigent: Enforcing QoS for latency-critical tasks on shared multicore systems," in *Proc. 21st Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2016, pp. 33–47.
- [68] D. K. Tam, R. Azimi, L. B. Soares, and M. Stumm, "RapidMRC: Approximating L2 miss rate curves on commodity systems for online optimizations," in *Proc. 14th Int. Conf. Architect. Support Program. Lang. Operating Syst.*, 2009, pp. 121–132.
- [69] X. Zhang, S. Dworkadas, and K. Shen, "Towards practical page coloring-based multicore cache management," in *Proc. 4th Eur. Conf. Comput. Syst.*, 2009, pp. 89–102.
- [70] N. Beckmann and D. Sanchez, "Jigsaw: Scalable software-defined caches," in *Proc. 22nd Int. Conf. Parallel Architect. Compilation Techn.*, 2013, pp. 213–224.
- [71] X. Hu, X. Wang, Y. Li, Y. Luo, C. Ding, and Z. Wang, "Optimal symbiosis and fair scheduling in shared cache," *IEEE Trans. Parallel Distrib. Syst.*, vol. 28, no. 4, pp. 1134–1148, Apr. 2017.
- [72] Y. Xiang, X. Wang, Z. Huang, Z. Wang, Y. Luo, and Z. Wang, "DCAPS: Dynamic cache allocation with partial sharing," in *Proc. 13th Eur. Conf. Comput. Syst.*, 2018, pp. 1–15.
- [73] O. A. Ben-Yehuda, E. Posener, M. Ben-Yehuda, A. Schuster, and A. Mu'alem, "Ginseng: Market-driven memory allocation," in *Proc. 10th Int. Conf. Virt. Execution Environ.*, 2014, pp. 41–52.
- [74] L. Funaro, O. A. Ben-Yehuda, and A. Schuster, "Ginseng: Market-driven LLC allocation," in *Proc. USENIX Annu. Tech. Conf.*, 2016, pp. 295–308.



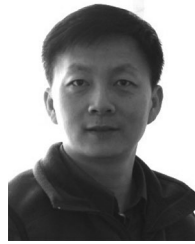
Taekyung Heo received the BS degree in computer engineering from Sungkyunkwan University, Seoul, South Korea, and the MS degree in computer science from KAIST, Daejeon, South Korea. He is currently working toward the PhD degree with the School of Computing, KAIST, Daejeon, South Korea. His research interests include memory systems, computer architecture, and accelerators.



Yang Wang (Student Member, IEEE) received the BS degree from the University of Electronic Science and Technology of China, Chengdu, China. He is currently an internship with Microsoft Research Asia. His research interests include computer architecture and general-purpose graphics processing unit architecture.



Wei Cui received the BS degree from the Nanjing University of Science and Technology, Nanjing, China, and the MS degree from Peking University, Beijing, China. He is currently a senior research SDE with Microsoft Research, Asia, Beijing. His research interests include computing accelerators, AI platform, and system optimization.



Lintao Zhang (Senior Member, IEEE) received the BS degree in physics from Peking University, Beijing, China, and the PhD degree in computer engineering from Princeton University, Princeton, New Jersey. He is currently a research manager with Microsoft Research Asia. His research interests include system issues in very large-scale, distributed systems.



Jaehyuk Huh (Member, IEEE) received the BS degree in computer science from Seoul National University, Seoul, South Korea, and the MS and PhD degrees in computer science from the University of Texas, Austin, Texas. He is currently a professor with the School of Computing, KAIST. His research interests include computer architecture, parallel computing, virtualization, and system security.

▷ **For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.**