# From Good to Great: Improving Memory Tiering Performance Through Parameter Tuning

Konstantinos Kanellis*
kkanellis@cs.wisc.edu
University of Wisconsin-Madison
Madison, WI, USA

Sujay Yadalam*
sujayyadalam@cs.wisc.edu
University of Wisconsin-Madison
Madison, WI, USA

Fanchao Chen
fcchen@cs.wisc.edu
University of Wisconsin-Madison
Madison, WI, USA

Michael Swift
swift@cs.wisc.edu
University of Wisconsin-Madison
Madison, WI, USA

Shivaram Venkataraman
shivaram@cs.wisc.edu
University of Wisconsin-Madison
Madison, WI, USA

## Abstract

Memory tiering systems achieve memory scaling by adding multiple tiers of memory wherein different tiers have different access latencies and bandwidth. For maximum performance, frequently accessed (hot) data must be placed close to the host in faster tiers and infrequently accessed (cold) data can be placed in farther slower memory tiers. Existing tiering solutions employ heuristics and pre-configured thresholds to make data placement and migration decisions. Unfortunately, these systems fail to adapt to different workloads and the underlying hardware, so perform sub-optimally.

In this paper, we improve performance of memory tiering by using application behavior knowledge to set various parameters (*knobs*) in existing tiering systems. To do so, we leverage Bayesian Optimization to discover the good performing configurations that capture the application behavior and the underlying hardware characteristics. We find that Bayesian Optimization is able to learn workload behaviors and set the parameter values that result in good performance. We evaluate this approach with existing tiering systems, Hemem and HMSDK. Our evaluation reveals that configuring the parameter values correctly can improve performance by 2x over the same systems with default configurations and 1.56x over state-of-the-art tiering system.

## 1 Introduction

**Need for memory tiering.** Modern data-intensive applications like graph processing, machine learning and in-memory databases demand large amounts of memory for high performance. Due to the scaling limitation of DRAM, this demand has lead to memory becoming one of the most significant costs of datacenters [30, 43]. One way to solve this problem is to supplement existing DRAM with memory tiers consisting of new memory technologies such as Non-Volatile Memory (NVM) [8] and CXL-based memory [19]. Memory tiering enables building systems with vast amounts of memory in a cost-effective manner.

Newer memory technologies such as NVM or CXL memory have different characteristics (speed, size, cost) compared to DRAM. The higher access latencies and lower bandwidths offered by these technologies can thus significantly degrade application performance. To alleviate this problem, there is a need to design memory tiering systems that make smart data placement decisions. They should keep frequently accessed (hot) data in faster tiers like DRAM, and infrequently accessed (cold) data in slower tiers, such as NVM or CXL.

**Limitations of existing systems.** There have been numerous works that aim to smartly place and migrate data in memory tiering systems transparently to the running applications. Unfortunately, these solutions fail to perform well under all circumstances; they either do not work well with certain applications or on certain system configurations. For example, prior work [24] has observed that HeMem [35] fails to place all the hot pages in fast tier for some workloads (also discussed in Section 4.2). The are two main reasons why these systems behave sub-optimally.

*First*, these solutions use heuristics, such as access frequency or recency of access, along with pre-configured thresholds to make data placement decisions. These heuristics and thresholds are optimized for the common case and not for specific workload behavior. As a result, these systems fail to correctly identify the hot pages and/or fail to migrate them in time. Memtis [24] attempts to address this limitation by dynamically adjusting the threshold for page promotion to ensure the correct set of hot pages are loaded in the fast tier. However, we find that even Memtis makes sub-optimal decisions in some cases (e.g. fails to identify write-heavy pages, see Section 4.6).

*Second*, the policies fail to capture and/or utilize high-level information about the application such as whether the access pattern for a region is sequential or random. For instance, HeMem migrates some write heavy pages during the initialization phase of an application which turn out to be un-useful during the more critical read phase of the application. If HeMem has knowledge about the application behavior, it could avoid migrating those pages.

---

*Both authors contributed equally to this work.

We believe that a memory tiering system that can adapt to the workload running and the underlying hardware characteristics can achieve better performance than existing solutions. In this paper, we try to answer the question: *how much performance improvement can be achieved by taking into account application behavior while making data placement decision in memory tiering systems?*

To this end, we test this hypothesis by tuning parameters, a.k.a *knobs*, of existing memory tiering systems. Our key insight is that different data placement and migration patterns can be realized by tuning various parameters in existing tiering systems. As an example, consider HeMem [35], which has a parameter called the `read_hot_threshold` which is the minimum number of accesses required to a page before considering the page a candidate for promotion from slow tier to fast tier. HeMem has another parameter called `migration_period` which determines when the candidate pages for promotion or demotion are migrated between memory tiers. Using these two parameters, one can control which pages are considered hot, and when they are migrated. Therefore, by adjusting knob values, one could possibly realize a memory tiering system that imitates and performs close to an optimal data placement algorithm [49].

Manually finding the best values for the the knobs is, however, extremely challenging. Not only is it time-consuming but even experts with domain expertise could fail to identify the relationship between the application behavior, the desired memory tiering system behavior, and its parameters. So, we leverage *Bayesian Optimization* to identify the best-performing tiering parameter configurations for different workloads. Bayesian optimization (BO) is effective at finding the optimal value of a function with few samples. In the past, researchers have employed BO to determine the best set of knob values in many systems such as databases [22] and storage systems [5]. Our work is the first to evaluate BO to tune knobs for memory tiering systems.

We use the optimizer to discover the best knob values for two existing tiering systems, HeMem [35] and HMSDK [23] under different scenarios by changing workloads, input datasets, thread counts etc. From our evaluation, we find that the optimizer is able to identify memory access patterns of applications and adjust knob values accordingly. For instance, the optimizer recognizes workloads with streaming access patterns and generates a HeMem configuration that avoids migrations of such non-beneficial pages. In comparison, HeMem with the default configuration keeps migrating pages leading to lower performance (see Table 5. Overall, we find that using an optimizer to tune tiering system parameters can yield as much as 2x improvement.

Finally, we analyze the best-performing configurations across scenarios using memory access patterns and migrations seen over time. By comparing these patterns with the default configuration we derive some common reasons why existing tiering systems fall short. We summarize these reasons in Section 5 and describe some research directions which can lead to better tiering systems in the future.

## 2 Background and Motivation

***Memory tiering.*** Modern applications such as in-memory databases, web serving, graph processing and machine learning demand large memory systems for high performance. They require high bandwidth and often have low tail latency requirements. Unfortunately, scaling memory naively has become challenging because: (1) DRAM scaling has stagnated over the past few years and the cost of DRAM has been increasing making memory costs a large portion of the Total Cost of Ownership (TCO) [43], (2) the number of channels per socket is limited by number of pins on the chip.

One way to achieve cost-effective memory scaling is through *memory tiering*, which involves the organization of data across multiple *tiers* of memory. New memory tiers are built using new memory technologies such as NVM [8] (lower cost/byte), CXL-based memories [19] or byte-addressable NVMe SSDs [45]. Unlike processor caches, in tiering, data resides in a single location (tier) in the memory hierarchy and is accessed directly from the tier it is present in. These new memory tiers exhibit higher latencies and/or lower bandwidth than locally connected DRAM [20, 40]. Applications can suffer from performance degradation if they access these slower tiers frequently. Therefore, for improved performance, it is desirable that the majority of memory accesses are made to data in faster tiers. To achieve this, memory tiering systems aim to smartly place frequently accessed hot data in tiers closer to the CPU and less frequently accessed cold data in slower farther tiers.
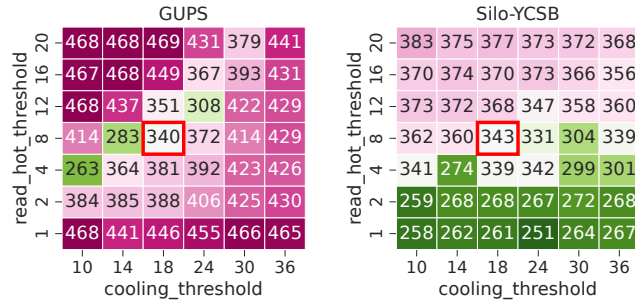
Researchers have proposed several approaches to smartly place data in memory tiers for achieving good performance [3, 13, 24, 28, 35, 44, 46, 47]. These systems perform two key tasks: (1) *identify hot and cold data/pages*, and (2) *migrate the identified hot/cold pages between tiers*. Ideally, for maximum performance, these systems should imitate data placement and migration that an ideal tiering algorithm such as $CH_{opt}$ [49] would perform.

***Limitations of existing tiering systems.*** Existing systems, unfortunately, fall short for two main reasons: (1) they employ heuristics and static thresholds as shown in Table 1 to identify and migrate hot/cold pages across tiers, and (2) they rely on low-level signals such as hardware event sampling (e.g. Intel PEBS [7]) or page table accessed (A) bit scanning (e.g. DAMON [32]) or page faults which do not provide high-level information such as the access pattern of the workload.

***Case study.*** To understand the extent to which heuristics and especially static thresholds could affect memory tiering performance, we perform a simple case study using HeMem [35] on an NVM machine. HeMem and other tiering systems often include several parameters, a.k.a *knobs*,

| | Access monitoring | Promotion heuristic | Promotion threshold | Demotion heuristic | Demotion threshold | Migration heuristics |
|---|---|---|---|---|---|---|
| AutoTiering [47] | Page fault | Access recency | Static threshold | Access frequency | | Watermark (demotion) |
| Thermostat [3] | Page fault | Access frequency | Dynamic threshold | Access frequency | Dynamic threshold | Migration period |
| AMP [17] | PT scanning | Recency/frequency/random | Static threshold | Recency/frequency/random | Static threshold | Migration period |
| Multi-clock [27] | PT scanning | Recency+frequency | Static threshold | Access recency | Static period | Migration period (promotion) Watermark (demotion) |
| HeMem [35] | PEBS | Access frequency | Static threshold | Access frequency | Static threshold | Migration period |
| TPP [28] | Page fault | Recency | Static threshold | Access recency | Static period | Watermark (demotion) |
| TMTS [13] | PT scanning+PEBS | Recency+frequency | Static threshold | Access recency | Static period | Migration period |
| HMSDK [23] | PT scanning (DAMON) | Access frequency | Static threshold | Access frequency | Static threshold | Migration period |
| Memtis [24] | PEBS | Access frequency | Dynamic threshold | Access frequency | Dynamic threshold | Migration period |

**Table 1.** Summary of heuristics and thresholds used in existing tiering systems. We list only the most important parameters, however each system has a lot of more parameters such as scanning/sampling frequency, page size granularity, size of lists, etc.



**Figure 1.** Execution time (in seconds) of *GUPS* (left) and *Silo* (right) workloads, when we tweak two HeMem parameters. Default configuration execution time is shown in red box.

such as hotness thresholds and migration period which control the system behavior (see table 2). Developers of these systems perform limited sensitivity studies to find the best values for these parameters which we refer to as the default values. We consider just two of HeMem's configuration parameters: `read_hot_threshold`, which is used to determine when a page is hot (default is 8 sampled accesses), and `cooling_threshold`, which controls when page counts are cooled, i.e., when page counts are halved so as to imitate an exponential moving average (default is 18 sampled accesses). We perform a straightforward grid search, where we tweak the values of these knobs to explore how larger and smaller values affect the execution time of a given workload; we use the default value for the remaining knobs. Figure 1 shows the execution time for GUPS [33] and Silo [42] (running a read-only YCSB-C workload) under different HeMem parameter configurations. We highlight the default HeMem configuration with a red box.

From our results, we make the following observations. First, we observe that different HeMem configurations result in large variations in workload performance (i.e., elapsed time). Second, we observe that there exist configurations that perform much better compared to the default HeMem

configuration for a given workload. For example, the best-performing configuration delivers up to 29% performance improvement compared to the default for GUPS and the improvement is 36% for Silo.

> **Insight #1** *Memory tiering system parameters such as thresholds and migration periods are critical for performance. Setting these parameters to the appropriate values can yield significant performance gains.*

Third, we observe that the knob values that improve performance are quite different across our workloads, GUPS and Silo. For GUPS, better performance is achieved with a large read hot threshold ($12 - 20$), while for Silo performance is improved when using small read hot threshold ($1 - 2$). Existing solutions such as HeMem have a myopic view wherein they fail to capture the high-level workload behavior and the characteristics of the underlying hardware. Thus, these systems fail to adapt to the environment they are running in and behave the same way under all scenarios which leads to sub-optimal decisions.

> **Insight #2** *The best values of these parameters depend on the workload. Since these parameters control the data placement and migration decisions made by the memory tiering system, they can be modified to express the desired migration sequence/pattern for an application.*

Memtis [24] makes a similar observation that static thresholds do not reflect the workload behavior and proposes a strategy to dynamically decides the thresholds for hot, warm and cold pages based on the workload access distribution. This allows Memtis to outperform previously proposed tiering systems. However, we find that dynamically adjusting only the hot/warm/cold thresholds is insufficient in many cases. Memtis also fails to consider the characteristics of the underlying system. The cost of page migrations can vary

depending on the available bandwidth and volume of competing traffic. Since Memtis does not consider this cost while making decisions, we see that it ends up making disadvantageous migrations. Also, Memtis uses static values for various other parameters such as the `cooling_threshold` and `migration_period` which leads to poor performance under many circumstances. (more details in Section 4.6)

Based on these observations we see two key challenges for memory tiering systems: first, it is important to efficiently *discover good performing parameter configurations* for existing tiering systems that can capture the application behavior and the underlying hardware characteristics. Second, it is crucial to *analyze why certain parameter configurations perform well*. Analyzing the best performing configuration can help us understand workload-specific migration patterns which lead to good performance and yield guidelines which can influence the design of future memory tiering systems.

## 3 Tuning Tiering System Parameters

As discussed in the previous section, finding the right set of values for tiering system parameters can result in major performance improvements. Prior works [13, 35] use domain expertise and sensitivity studies to find the "best" parameter values but this does not scale well. Additionally, manually finding the best values for every workload and hardware platform can be extremely challenging.

Some tiering systems also have a large parameter space (e.g., many parameters with large value ranges), with several parameters being correlated with each other [23, 35]. Therefore it is possible that only a specific combination of right parameter values may actually bring performance benefit, meaning that parameters should not be tuned independently from each other. Moreover, there might exist *hidden* parameters, i.e., parameters that are not explicitly exposed by the tiering engine. Such parameters are often neglected, but if considered, may actually bring meaningful improvements (as we show for Silo in Section 4.2). Note that the above challenges are not limited to memory tiering engines, but have also been observed in many other systems, including DBMSs [15, 22] and big-data analytics frameworks, like Spark [14], or key-value stores, like RocksDB [29, 50].

Furthermore, understanding the relationship between application behavior, parameter values and their effect on the migration decisions can be quite tricky. To address this, we need a tuning process that is more principled compared to traditional grid- or random- search methods, which explore the parameter space in a rather unguided fashion. Instead, we need a process that (1) *maximizes* the information gained from each configuration evaluation, (2) is very *sample-efficient* and (3) continuously tries to *learn* how the tiering engine behaves.

### 3.1 Proposal: Leverage Bayesian Optimization

To this end, we propose employing Bayesian Optimization (BO). BO is a sequential, model-based optimization process that aims to find the (global) optimum of a black-box function using as few trials as possible [39]. It is typically used when the black-box function is expensive to evaluate. Notably, BO does not make any a-priori assumptions regarding the shape/behavior of the function (e.g., convex, monotonic, etc), making it a prime candidate to use in our case [37].

In our work, BO can be used to identify the best-performing configuration (global optimum) for our desired tiering engine (black-box function). To find a best-performing configuration, we need to determine an appropriate value for each parameter based on the workload and hardware. Formally, given a set of tiering engine parameters $\theta_1, ..., \theta_n$ along with their domains (i.e., range of values), $\Theta_1, ..., \Theta_n$, the tiering system parameter space can be defined as $\Theta = \Theta_1 \times ... \times \Theta_n$. We want to find a configuration $\theta^* \in \Theta$ that minimizes the execution time of the workload $f$:

$$\theta^* = \arg\min_{\theta \in \Theta} f(\theta)$$

A Bayesian optimizer consists of two core components: the surrogate model and the acquisition function. The surrogate model is a machine-learning model that models the behavior of the system (i.e., tiering system), and can predict its performance given a specific configuration. Initially the model has no knowledge of the system, but with more and more observations (i.e., evaluated configurations), it quickly "learns" how the tiering engine behaves. The acquisition function is used by the optimizer to decide which configuration is more promising to evaluate next. This is done by balancing exploration (i.e., evaluate a configuration on a previously unexplored region of the parameter space) with exploitation (i.e., try to improve an already good-performing configuration by making small tweaks on parameter values). **Tuning Pipeline.** Given a fixed workload, we first launch a tuning session, where we (1) set our tiering sytem to a configuration (initially default), (2) execute the workload, and (3) measure its execution time. Then, we feed the evaluated configuration alongside the execution time back to the optimizer. The optimizer updates the surrogate model with the new observation and selects the best candidate configuration to be evaluated next (via the acquisition function). We repeat the above process, always using the new configuration suggested by the optimizer, until we exhaust the tuning budget (e.g., number of iterations). To balance between exploration and exploitation, we also (1) initially bootstrap the optimizer with some randomly-sampled configurations, and (2) force the optimizer to periodically suggest a random configuration (i.e., random config probability).
**BO Algorithm.** In this work we use the state-of-the-art Sequential Model-based Algorithm Configuration (SMAC)

| Knob Name | Default | Min | Max | Description |
|---|---|---|---|---|
| sampling_period | 5000 | 100 | 10000 | Number of memory load events to trigger sampling |
| write_sampling_period | 10000 | 1000 | 20000 | Number of store instructions to trigger sampling |
| read_hot_threshold | 8 | 1 | 30 | Minimum number of read access samples per page to classify it hot |
| write_hot_threshold | 4 | 1 | 30 | Minimum number of write samples per page to classify it hot |
| cooling_threshold | 18 | 4 | 40 | Number of sampled accesses to trigger page access count cooling |
| migration_period | 10 | 10 | 5000 | Interval of migration thread executions (ms) |
| max_migration_rate | 10 | 2 | 20 | Maximum migration rate allowed (GiB/s) |
| cooling_pages | 8192 | 1024 | 65536 | Number of pages cooled at a time |
| hot_ring_reqs_threshold | 1024 | 128 | 4096 | Number of hot pages processed at a time |
| cold_ring_reqs_threshold | 32 | 8 | 256 | Number of cold pages processed at a time |

**Table 2.** HeMem parameter knobs, their defaults, and ranges input to the optimizer

framework [18], which is a popular BO-based optimizer [37]. SMAC utilizes a Random Forest (RF) as the surrogate to model the parameter space; this enables it to handle large high-dimensional spaces efficiently. As shown by prior work [22], when used to tune complex systems, SMAC manages to find best-performing configurations within few evaluations, even when having hundreds of knobs, with a subset of them correlated with each other. Considering our previous case study in Figure 1, SMAC was able to find the best-performing configuration for GUPS within 10 − 16 iterations, making it 2.5 − 4× more sample-efficient.

An additional benefit of using SMAC's random forest surrogate model is its inherent ability to determine which tiering system knob(s) are more important. This is performed by computing an *importance score* for every knob, based on its impact to the workload execution time (similar to [5, 21]). In particular, for each knob $k$, we fix the values of other knobs at their default value and then randomly sample across values of $k$ to observe the impact on workload performance (via surrogate model prediction). Assuming enough observations, the surrogate model can make accurate predictions. While one should not blindly rely on these importance scores, these can be helpful when trying to understand why some configurations perform better than others, like in Section 4.2.

### 3.2 Tiering engine under the spotlight: HeMem

We use the proposed optimization pipeline to tune knobs of two tiering engines, HeMem [35] and HMSDK [23]. In this section, we give a brief overview of HeMem which will be helpful to understand our analysis in Section 4. We discuss HMSDK later in Section 4.5. HeMem is a user-space memory manager for tiered memory systems that is dynamically and transparently linked into applications. HeMem intercepts mmap(), handling page allocations pertaining to anonymous memory. Table 2 lists all configuration knobs we consider, which are also discussed in the following paragraphs.

HeMem monitors page accesses using hardware event sampling such as Intel PEBS [7]. HeMem samples L3 load misses as well as all store instructions. Based on the samples gathered, HeMem maintains an access count per page. HeMem uses separate counters for reads and writes. If the page access count exceeds a pre-determined threshold, it classifies the page as hot; else considers the page cold. It uses different thresholds for reads (read_hot_threshold) and writes (write_hot_threshold).

HeMem uses a page count cooling mechanism to give importance to latest accesses. When the access count of any page reaches a cooling_threshold, page cooling is triggered which halves access count of all pages. If the access count falls below the hot threshold after cooling, the page is then considered cold. To avoid scanning all pages, HeMem cools pages in batches whose size is determined by another parameter called cooling_pages. This is one of the (hidden) parameters that is not discussed in the paper but is part of the implementation.

HeMem includes a background migration thread that executes periodically every migration_period. The migration thread promotes hot pages on slower tiers and demotes cold pages from faster tiers. HeMem marks pages as write-protected when they are being migrated to avoid races between application writes and the data movement. Applications could therefore stall waiting for pages to be migrated. **Deployment Issues.** During our initial tests with HeMem, we observed some abnormal behavior, which often led to workload crashes or poor performance, affecting our results. In an effort to achieve a more accurate evaluation, we tried to address these issues without compromising the original HeMem implementation design and functionality. Below, we briefly summarize these issues and our corresponding fixes:

1. *High sampling overhead*: In the original code, PEBS hardware generated an interrupt for every sample, incurring high overhead. This is due to certain parameters requested from PEBS (PERF_SAMPLE_WEIGHT). We modified this parameter, so that the hardware generates much fewer interrupts without affecting sampling accuracy.
2. *Migrate-free page race condition*: HeMem's original implementation does not properly protect against cases where a page is freed by the workload, while HeMem tries to migrate at the same time. This causes some workloads to completely crash. We fixed this issue by introducing lightweight synchronization.

| Specification | pmem-large | pmem-small | NUMA |
|---|---|---|---|
| Number of cores | 24 | 16 | 20 |
| Processor generation | Icelake | Cascadelake | Skylake |
| Processor frequency (GHz) | 3 | 3.2 | 2.2 |
| L3 cache size (MB) | 18 | 11 | 13.75 |
| Far memory type | Optane | Optane | NUMA |
| Max near memory size (GB) | 96 | 32 | 96 |
| Max far memory size (GB) | 128 | 128 | 96 |
| Max near mem BW (GB/s) | 138 | 46 | 56 |
| Max far mem BW (GB/s) | 7.45/2.25 | 6.8/1.85 | 36/36 |
| Near memory latency (ns) | 80 | 80 | 95 |
| Far memory latency (ns) | 150 - 250 | 150 - 250 | 145 |

**Table 3.** Specifications of machines used in our evaluation.

3. *Large minimum allocation size*: By default, HeMem is configured to only handle page allocation greater than 1GiB in size. We changed this to 128MiB, as we observed that some of our workloads do not make allocations that large, making HeMem redundant for those cases.

4. *Read/write sampling frequency separation*: We introduce a new parameter called WRITE_SAMPLING_PERIOD that controls PEBS sampling frequency of store instructions. This allows HeMem to exert different importance for read and write operations.

## 4 Evaluation & Analysis

Through our evaluation and analysis, we try to answer the following questions:

1. How much performance improvement can the best-performing configuration achieve for HeMem across workloads?

2. What are the important differences between the default and the best-performing knob values? How does this translate into additional performance?

3. Is there a single best-performing configuration that would work well in all scenarios?

4. Under what circumstances do tiering systems leave unexploited performance on the table?

5. How do the above results and observations generalize to a different tiering system?

6. How does the performance of a tuned tiering system compare against state-of-the-art designs like Memtis which dynamically adjust certain knobs (e.g. hot_threshold)?

### 4.1 Experimental Setup

***Hardware setup.*** Table 3 shows the hardware specifications of the three machines used in our evaluation. pmem-large and pmem-small have Intel Optane DC Persistent Memory DIMMs which form the slower tier, whereas the NUMA machine (provided by Cloudlab [12]) is used for emulating CXL memory. Notably, pmem-small has 3× less DRAM bandwidth compared to pmem-large. Unless mentioned otherwise, we run our experiments on pmem-large.

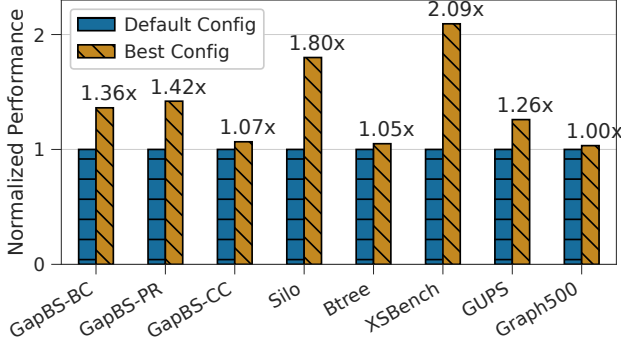| Workload | Inputs | RSS | Description |
|---|---|---|---|
| GapBS-BC [4] | kronecker twitter | 78.13 13.08 | Compute the measure of centrality in a graph based on shortest paths. |
| GapBS-PR [4] | kronecker twitter | 71.29 12.32 | Compute the PageRank score of a graph |
| GapBS-CC [4] | kronecker twitter | 69.29 12.09 | Compute connected components of a graph using (Shiloach-Vishkin) |
| Silo [42] | TPC-C YCSB-C | 75.68 71.40 | In-memory transactional database |
| Btree [1] | - | 12.13 | In-memory index lookup benchmark |
| XSBench [41] | - | 64.97 | Compute kernel of the MCNP algorithm |
| GUPS [33] | 8 GiB hot | 64.03 | Random accesses with dynamic hotset |
| Graph500 [31] | kronecker | 34.13 | Construction and BFS of large graphs |

**Table 4.** Workload Characteristics. RSS is in GiB.

***Workloads.*** For our evaluation, we select a set of 8 representative and diverse workloads. Specifically, we employ three graph processing algorithms from GapBS [4], Graph500 [31], an HPC workload (XSBench [41]), an in-memory database (Silo [42]), an in-memory index lookup (Btree [1]), and a popular memory intensive micro-benchmark, GUPS [33]. Table 4 contains additional information about these workloads, including the resident set size (RSS) and the different inputs used (if any) for each workload. This set of workloads have been widely used to evaluate tiered memory systems in many prior works [24, 25, 35, 47]. When measuring execution time, we make sure to only consider relevant applications phases (e.g., graph construction/traversal, query execution), ignoring any potential reading/initialization phases that are not part of the actual kernel (e.g., disk I/O). We run these workloads with a thread count large enough to just saturate the memory bandwidth of each system, i.e., 12 threads for pmem-large and NUMA, and 4 threads for pmem-small.

***Tiering Configuration.*** Similar to Memtis [24], we configure the ratio of fast to slow tier memory size by setting the fast tier size to the corresponding proportion of the workload resident set size (RSS). Unless otherwise noted, experiments maintain a 1:8 memory size ratio. For instance, for GUPS whose RSS is 64 GB, we set the fast tier size to 7.11 GB (11%). In HeMem, we set the fast tier size as a macro in a header file. In HMSDK, we run a background *memchew* workload that reserves the desired portion of the fast tier, but does not perform any memory operations.

***Optimizer Configuration.*** We configure SMAC to run with a budget of 100 iterations, using the first 20 as an initial exploration phase; we also force the optimizer to pick a random configuration with 20% probability. Similar setups have been used in prior works, for more complex systems [22]. HeMem exposes its knobs as macros in the library. At every iteration, in order to deploy the newly-suggested system configuration, the optimizer modifies the values of these macros and recompiles the library before starting the execution. HMSDK can be configured by passing a JSON file as input.

**Figure 2.** Performance improvement of best HeMem configuration found over default for all workloads.
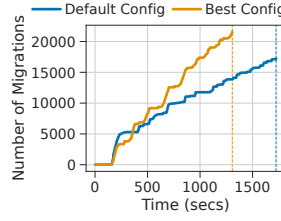


**Figure 3.** *GapBS-BC*: Number of migrations over time.



**Figure 4.** *GapBS-PR*: Memory access pattern over time.

| | Reason(s) for improvement | Important knobs |
|---|---|---|
| BC | Accurate classification of hot/cold pages, timely migrations | read_hot_threshold, sampling_period, cooling_threshold |
| PR,CC | Eliminate unnecessary migrations | read_hot_threshold, write_hot_threshold, cooling_threshold |
| Silo | Better hot pages identification, reduce migrations of warm pages, lower write sampling overhead | read_hot_threshold, sampling_period, cooling_pages |
| Btree | Decrease importance of write-heavy pages | write_hot_threshold, write_sampling_period |
| XSBench | Eliminate warm page migrations | read_hot_threshold, write_hot_threshold, cooling_threshold |
| GUPS | Increase sampling accuracy, faster detection of hot pages | read_hot_threshold, sampling_period write_hot_threshold, sampling_period |

**Table 5.** Summary of performance improvement from parameter tuning on pmem-large for different workloads

## 4.2 Performance benefit of parameter tuning

We begin by comparing the performance of the best knob configuration found by the optimizer against the performance with the default HeMem configuration. We consider each workload in our benchmark set and also analyze why the optimizer configuration is better by relating to the workload behavior. Figure 2 shows the results. For almost all workloads barring Graph500, the optimizer identifies parameter values that provide superior performance, by $1.07 - 2.09x$.

Below, we briefly discuss the differences in parameter values and the reasons for these improvements, while Table 5 summarizes our findings.

***GapBS-BC (kronecker graph):*** We run Betweenness Centrality (BC) from the GapBS suite with a synthetic kronecker graph with 1 million vertices. We find that the best configuration found by the optimizer outperforms the default by about 1.36x. The optimizer tunes the hotness thresholds and cooling threshold which results in a more accurate classification of hot and cold pages, and timely migrations. Figure 3 shows how the number of migrations increases over time for the default and best configurations. With the default configuration (read_hot_threshold of 8), HeMem takes longer to identify all the hot pages, even fails to identify some of them. With the best configuration (read_hot_threshold of 4 or 5), HeMem identifies most of the hot pages quickly at the beginning of every iteration and promotes them to fast tier (each step in the graph corresponds to the beginning of an iteration). We also find that tuning the cooling_threshold is important to avoid frequent cooling which could trigger demotions of actually hot pages.

***GapBS-PR, CC (kronecker graph):*** With PageRank (PR) and Connected Components (CC), we find that a small set of pages are hot while the rest of the pages are cold and accessed infrequently. As seen in Figure 4, PR exhibits a streaming memory access pattern for the cold pages. The optimizer learns that PR does not benefit from migrations of these cold pages as there is little to no reuse. On the other hand, the default configuration keeps migrating pages throughout the experiment which eats into the memory
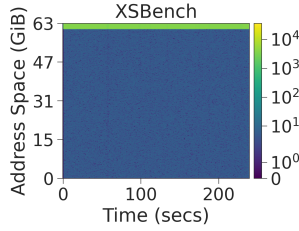
bandwidth and stalls many write accesses, slowing down the application. Similar to PR above, CC also has a streaming access pattern and doesn't benefit from migrations. The best-performing configuration is similar, in that it avoids any wasteful migrations.

> **Takeaway:** *Knowing the workload behavior and possibly the graph structure, in which vertices are popular while others are not, can help an ideal tiering system place appropriate pages in the faster tier and avoid all migrations.*
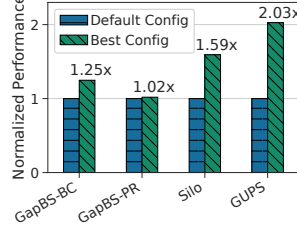
***Silo (YCSB-C):*** Silo is an in-memory transactional database. Similar to Memtis [24], we run Silo with YCSB-C read-only workload. With this workload, about 1% (700MB) of the pages of Silo are extremely hot, while about 20% of the pages are warm. The optimizer generates configurations that perform better than the default. Our analysis reveales three main reasons: (1) the best configuration identifies hot pages more accurately and quickly by reducing sampling_period and reducing read_hot_threshold), (2) it reduces migrations of moderately warm pages by cooling all pages at the same time rather than in batches (increases cooling_pages), and (3) reduces write sampling overheads, which makes sense since YCSB-C is a read-only workload.

The important knob analysis concluded that value of cooling_pages plays an important role in Silo's performance. This knob is one of the aforementioned hidden

**Figure 5.** *XSBench*: Memory access pattern over time.



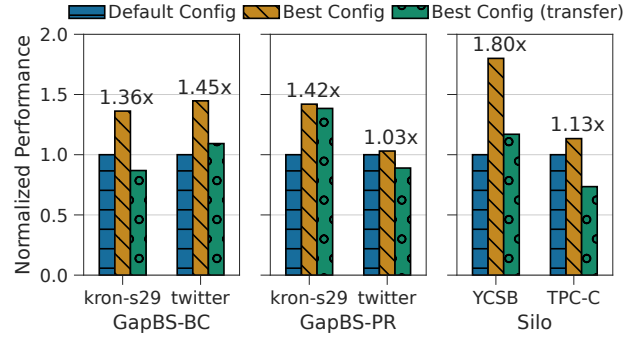**Figure 6.** Performance gains (vs default) for `pmem-small`.



**Figure 7.** Performance of HeMem's best configuration compared to default, for two application inputs. We also show how each best configuration performed when HeMem used it while running on the alternative input (i.e., *transfer*).

knobs that are usually neglected, highlighting the need to consider all possible knobs.

***Btree:*** We populate an in-memory Btree with 600M key-value pairs and perform 750M lookups. This workload has 2 phases: an initialization phase where keys are inserted into the Btree, and a lookup phase, where uniform random keys are queried. The initialization phase is write-heavy as new elements are created and added to the tree. Note that the tree might be adjusted after every insert to meet Btree properties. With the default configuration, HeMem starts migrating pages that are being written to. Since HeMem uses a low `write_hot_threshold`, a lot of pages get hot and are migrated. In fact, about $16,000$ migrations out of the $18,000$ migrations happen during this initialization phase. The optimizer realizes that the hot pages in the initialization phase will not remain hot during the lookup phase. Thus, it generates configurations with higher `write_hot_threshold`, so very few write-heavy pages are migrated during initialization. This creates room for read-hot pages such as those holding the high level nodes in the fast tier. This in turn reduces the latency of lookup queries.

> **Takeaway:** *Identifying that program phases are read-only (write-heavy) can help an ideal tiering system adjust the appropriate read hot (write hot) thresholds and minimize migrations.*

***XSBench:*** Figure 5 shows the access heatmap of XSBench Monte Carlo algorithm. XSBench has a small set of pages that are frequently accessed (greenish-yellow line at the top). Rest of the pages (blue area) are randomly accessed and have very similar access counts. Thus, it is important to keep the hot pages in the fast tier; it does not matter which of the remaining pages are also in the fast tier. The optimizer learns this and converges to configs that keep the hot pages in fast tier, avoiding wasteful migrations of the randomly accessed pages. To do so, it sets the `read_hot_threshold` and `write_hot_threshold` higher than the `cooling_threshold`, to not trigger any migrations.

***GUPS:*** We evaluate GUPS micro-benchmark with a skewed access distribution similar to prior work [25, 35]. We configure GUPS with a hotset of 8GiB and a total memory footprint of 64GiB. The hotset moves after half of the updates are performed. For our 1:8 memory configuration, the hotset spills

out of the fast tier. With the default configuration, HeMem keeps shuffling hot pages between the tiers because of low sampling accuracy. The optimizer recognizes this issue and increases the sampling frequency resulting in more accurate access sampling. This does not only improve hot page identification, but it also reduces the wasteful migrations opening up more bandwidth for the application.

> **Takeaway:** *Although increasing sampling frequency incurs higher overheads, some workloads can benefit more from the increased accuracy. So, constraining sampling overheads might not be a good idea for some applications.*
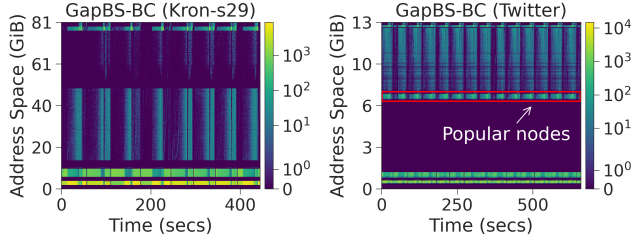
***Validating results on pmem-small:*** We run similar experiments to the above on our smaller pmem-small machine to verify that (1) performance gains are possible when switching to different hardware, and (2) the optimizer can identify similar best-performing configurations. Since pmem-small has smaller DRAM bandwidth, we run the same applications with fewer threads and smaller inputs. Figure 6 shows the performance difference between the default configuration of HeMem and the best configurations identified by the optimizer. Overall, the results are very similar to those obtained on pmem-large, corroborating the results discussed above.

### 4.3 Tuning with different application inputs

Application memory access patterns depend heavily on the input datasets/workload. For instance, the memory access pattern for graph processing applications could be very different for different graphs. Similarly, key-values (KV) stores exhibit different memory characteristics when running against different workloads (e.g. TPC-C versus YCSB-C). Since application behavior changes across inputs, the best-performing tiering knob values might also be different for different inputs. To answer this question, we ran our optimizer pipeline for graph and DB workloads with different inputs.

To understand if knob values are dependent on inputs, we use the best knob values obtained for one input with
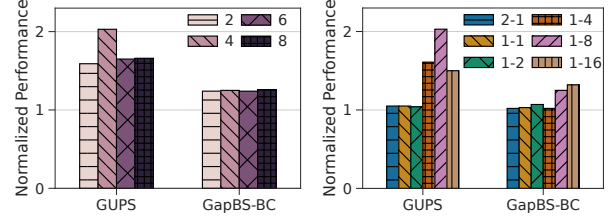
**Figure 8.** *GapBS-BC* memory access pattern over time, for Kronecker-s29 (left) and Twitter (right) input graphs.



(a) Varying Thread Count   (b) Varying Memory Ratios

**Figure 9.** Performance gains (over default) of HeMem's best configuration found, for varying (a) number of threads and (b) memory size ratios.

the other input. For example, we ran GapBS-BC on a twitter graph using the best configuration obtained for the kronecker graph. Figure 7 shows the results. As expected, in most cases, the best configuration generated for one input does not perform well for the other input. In fact, we observe that the performance is worse than default in most cases.

Below, we explain some of the key differences between the best configurations we found for the different inputs.

**Gap-BC, PR:** We compare the results of optimization runs for BC and PR with different inputs: kronecker and Twitter graphs. The memory access patterns is similar for both the graphs except for one difference (see Figure 8). With the twitter graph, there are a handful nodes corresponding to influential profiles that are very popular. These popular nodes are present on a small set of pages as shown in Figure 8. This results in frequent accesses to these pages whereas in kronecker, all pages receive almost equal number of accesses. The optimizer successfully identifies this difference. For the Twitter graph, the optimizer selects configuration which allow the the hot pages corresponding to the popular nodes to be migrated early. It does so by lowering `write_sampling_period` and the `write_hot_threshold` allowing for accurate and faster promotion of the popular pages.

**Silo:** We compare the best configurations for Silo with YCSB-C and TPC-C workloads. YCSB-C and TPC-C are very different workloads. While YCSB-C is a read-only workload with Zipfian access distribution, TPC-C is an insert-heavy workload that mostly reads new data, meaning that pages are hot when they are updated but become cold quickly as new entries are inserted into the database. As discussed earlier, for YCSB-C, the optimizer prioritizes read sampling parameters and de-prioritizes write sampling. On the other hand, for TPC-C, the optimizer recognizes the importance of writes. It arrives at configurations that promotes newly written pages to the fast tier for a short time and demotes them to slow tier as soon as they lose their hotness. The best-performing configurations (1) prioritize written pages by lowering write hot threshold and decreasing write sampling period, (2) invoke cooling frequently because hot set changes quickly, and (3) run the migration thread frequently by setting the migration period to the lowest value.

---

> **Takeaway:** *Application behavior is tied to the inputs they are processing. Thus, it is not sufficient to find just a single set of knob values for an application. System administrators need to consider the input properties along with the application when configuring knob values.*

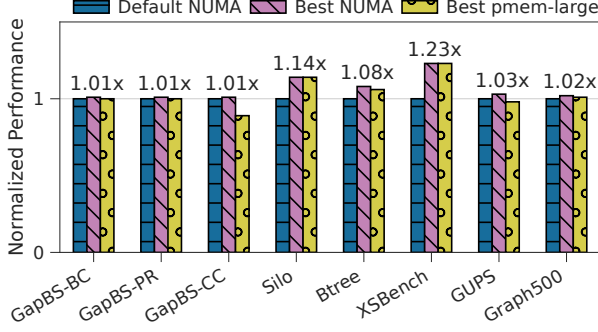### 4.4 Tuning for different system configurations

Memory tiering parameter values could depend on the properties of the underlying system. This includes various factors, significant among which are the number of cores, memory sizes and memory bandwidth. In this section, we look at the relationship between knob values and the aforementioned hardware characteristics.

**4.4.1 Different number of threads.** Programmers often configure the number of threads for an application based on the number of cores and the memory bandwidth available. Increasing the number of threads increases memory traffic which results in more accesses per page in shorter time period. Therefore, tiering engines need to have different thresholds and periods based on the number of threads.

Using two workloads, GUPS and BC (twitter graph), we perform optimization runs varying the number of threads. Figure 9a shows the performance improvement we achieve by changing knob values for different number of threads, on `pmem-small`. We see consistent performance improvement for all thread counts for both workloads. We see that the best knob values are *very different* for different thread counts. No single configuration performs well for all thread counts.

Through our analysis, we find that the `read_hot_threshold` is the most important knob and its value increases with number of threads, as expected. However, we find that it is not sufficient to just tune this single parameter when changing the number of threads. Other knobs such as the `migration_period` and `cooling_threshold` also have to be tuned accordingly.

**4.4.2 Different memory size ratios.** To understand the relationship between knob values and memory tiering size configurations, we run optimization on the knobs for different fast-slow tier memory size ratios. Recall that the memory

**Figure 10.** Performance gains (over default) of HeMem's *best* configuration for NUMA machine (annotated). We also show how `pmem-large` best config performed when deployed on NUMA machine.



**Figure 11.** Performance of best HMSDK configuration found (over default) for all workloads on NUMA machine.

NUMA machine (e.g. Silo, Btree, XSBench). It might therefore be possible to transfer the best-performing configuration from one machine to another machine for some workloads.

### 4.5 Tuning other memory tiering systems

As discussed in Section 2, almost all of the existing memory tiering systems use heuristics and have knobs of different kinds. In the previous sections, we saw how changing tiering parameter values improves HeMem's performance under various scenarios. To explore whether similar can be extended to other tiering systems, we also experiment with Heterogeneous Memory Software Development Kit (HMSDK) [23]. HMSDK is an open-source memory tiering system developed by SK-Hynix that uses DAMON [32] to identify hot and cold pages in memory. DAMON is a data access monitoring framework subsystem for the Linux kernel that scans page tables periodically. To reduce the overheads, DAMON divides the address space into a number of regions and samples page table entries from each region. HMSDK, like HeMem, has several knobs related to access monitoring, hot/cold page classification and migrations.

We setup HMSDK on the NUMA machine emulating CXL and run all workloads with 12 threads. Figure 11 shows the performance improvement achieved by changing various knob values in HMSDK. Again, we observe significant gains for some workloads and modest gains with others. Below, we discuss some key differences between the default knob values and the best knob values that result in better performance.

*GapBS-PR and Btree:* With the default knob values, HMSDK is unable to correctly identify hot pages and promote them. The optimizer improves page access monitoring by increasing both the number of DAMON regions (`nr_regions`) and the scanning frequency (`intervals:sample_us`). This leads to migrations inside the correct hot regions and reduces erroneous or unnecessary migrations, thereby improving performance.

*XSBench:* As discussed previously, XSBench has a small set of hot pages which should be placed in near memory while the rest of the pages have very similar access frequency. Similar to HeMem, we observe that HMSDK by default performs
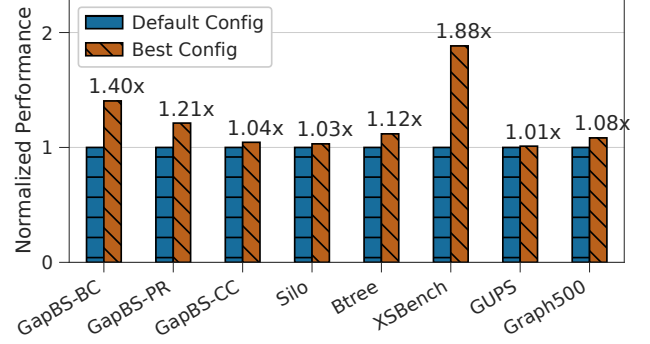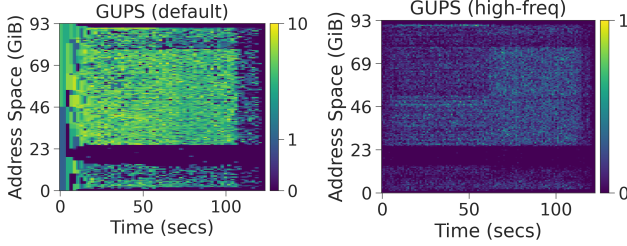
sizes of the fast and slow tier are configured based on the measured RSS of the workload.

Figure 9b shows the performance improvement over default knob setting for different memory ratios on `pmem-small`. The key observation here is that the HeMem configuration becomes more important for smaller fast tier sizes, where the entire hot set might not be able to fit in fast tier (e.g., 1:16, 1:8). Indeed, we find that the default configuration performs poorly in such settings. The optimizer is able to identify knob values that result in more accurate page classification and timely migrations. For instance, the optimizer sets the hot thresholds higher in the 1:16 and 1:8 tiering configurations and sets them lower in the 1:2, 1:1 and 2:1 configurations. As a result, only the extremely hot pages are migrated in the 1:16 and 1:8 systems whereas any page sampled/accessed is promoted in the larger fast tier configurations.

#### 4.4.3 Different memory bandwidth.
Studies on CXL show that CXL memories can offer latency and bandwidth comparable to NUMA [40]. Therefore, we emulate CXL memory using a NUMA machine wherein the local node memory represents fast tier and remote node memory emulates the slow tier. Table 3 shows the latency and bandwidth of local and remote NUMA memories.

Figure 10 shows the performance improvement achievable by tuning knobs of HeMem while running on the NUMA machine. We observe that the performance gains are mostly modest. The explanation is straight-forward: since latencies and bandwidth are similar between the tiers, there is little room for performance improvement. Second, migrations have little to no interference with application accesses because of high available bandwidth. Therefore, the penalty for potential wasteful migrations is negligible, so tiering systems do not need to be systematically tuned to perform well.

Interestingly, we also observe that best performing configurations on pmem-large mostly perform well even on the

**Figure 12.** Heatmap of GUPS generated by DAMON, using default HMSDK scanning frequency (left) and high scanning frequency (right)

a large number of unnecessary migrations (about 10 million pages). The optimizer identifies knob settings that eliminate migrations of such pages with similar access frequency, which improves performance significantly.

**GUPS:** We observe that it is not possible to improve HMSDK performance, which is surprising given that GUPS is a very simple workload. After investigating further, we found that HMSDK is unable to differentiate between heavily accessed pages and lightly accessed pages irrespective of its configuration. This constitutes a major limitation of DAMON, which splits the address space into regions and assumes that all pages within a region have similar access frequency. Unfortunately, the hot pages of GUPS are distributed across the address space. Thus, DAMON finds all the regions to be hot. Figure 12 shows the heatmaps generated by DAMON for GUPS with two different sets of monitoring parameters. Irrespective of the parameter values, DAMON is unable to identify hot and cold pages accurately.
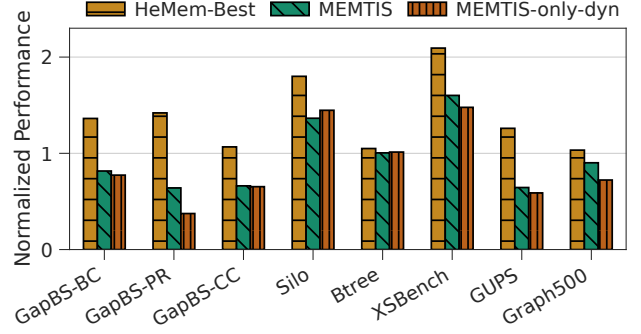
> **Takeaway:** *HMSDK experiences issues similar to HeMem, despite their apparent dissimilar design. HMSDK issue stem from static thresholds, making it impossible to reflect the application behavior.*

### 4.6 Comparison to Memtis

Memtis [24] is a state-of-the-art memory tiering system that improves upon HeMem. Memtis also employs PEBS to monitor page accesses and runs migration periodically in the background. Memtis' core improvements over HeMem are:

1. It *dynamically adapts hot thresholds* so as to maintain the hot set size close to the fast tier capacity.
2. Introduces a *warm* class for page classification. Memtis avoids migrating warm pages when the migration overhead would overshadow the benefit.
3. Memtis dynamically chooses the page size (hugepage versus regular page) based on subpage access skewness. If Memtis finds that only a small subset of pages in a hugepage are hot, it splits the huge-pages into base pages and promotes only the hot base pages.

We compare Memtis' performance with the one achieved by default and best HeMem configurations. We include



**Figure 13.** Performance of HeMem's best config found vs Memtis, normalized to HeMem default config performance.

Memtis' performance, where we enable only the dynamically threshold adaptation feature and disable the warm pages and hugepage split features (i.e., MEMTIS-only-dyn). Figure 13 shows the normalized performance for our workload set.

Similar to the results shown by authors of Memtis, we see that Memtis outperforms HeMem on default configuration for some workloads (e.g., Silo, XSBench), yet for most Memtis is slower. More importantly, we observe that the best-performing HeMem configuration outperforms Memtis for all workloads.

Our analysis revealed multiple factors for the poor performance of Memtis. We find that Memtis spends a significant amount of time in the kernel for page allocations, page splitting and migrations for these workloads. Second, Memtis uses very high sampling period (100K) for writes which leads to low write sampling accuracy. Therefore, Memtis performs poorly on workloads that are write-intensive. Third, Memtis uses static thresholds for other parameters such as cooling period, hot threshold adaptation period and migration period. By only dynamically adjusting the hot threshold, Memtis is unable to fully adapt to the workload.

> **Takeaway:** *Dynamically adjusting one or two knobs is beneficial but insufficient to extract maximum performance from tiering systems. For maximum performance, tiering systems need to adjust all parameters to reflect the application behavior.*

## 5 Discussion and future work

One of the key takeways from our work is that existing tiering solutions do not perform well under all circumstances. They fail to adapt to the workload and/or the underlying hardware. Going forward, we need to build tiering systems that can use high-level application knowledge and monitor hardware characteristics while making data placement decisions. Ideally, we should design memory tiering systems that are fully adaptable, i.e., have no knobs. However, designing such a complex system can be quite challenging. Below, we list some ways to get close to such a solution:

**Use programmer annotations/hints.** One of the main takeaways from our evaluation is that utilizing workload behavior to make data placement decisions in memory tiering systems can provide significant benefits. We realize this by using an optimizer which learns the application behavior and sets tiering parameters accordingly. However, using an optimizer has some limitations (e.g., it is time consuming, requires access to the application and inputs) and is not feasible to use in all scenarios. An easier way to provide high-level information to the tiering engine is through program annotation and hints. In some of the workloads we evaluated such as Graph500 and XSBench, we find that the hot and cold pages are from different allocations. Programmers can therefore provide hints at allocation time about the type of access: random or sequential and about the expected hotness which would help tiering systems make smart placement and migration decisions.

**Use cost-benefit models.** Prior work [26] has argued that the kernel should use cost-benefit models to make policy decisions. We believe this argument holds for memory tiering systems as well. Through our evaluation, we find that existing systems [23, 24, 35] make sub-optimal policy decisions because they ignore the cost of page migrations. For many workloads (e.g. PageRank and XSBench), HeMem with the default configuration performs many high-cost non-beneficial page migrations. Migrations are expensive on machines with limited memory bandwidth or when applications are saturating the bandwidth. Page migrations, especially under memory pressure, has a detrimental impact on application performance [44]. Tiering systems should estimate the cost of page migration and benefit from the low latency access before moving pages between tiers. One way to estimate the cost is to measure the bandwidth utilization or memory queuing latencies using hardware performance counters.

**Identifying and adapting to application phases.** Many workloads have distinct execution phases. For instance, machine learning algorithms have a data loading phase, a preparation phase, and then in every iteration a forward propagation and back propagation. The memory access behavior is very different in each of these phases. Existing tiering systems are unable to identify and adapt to the different phases. For instance, with the Btree workload, we find that the optimizer fails to accelerate both the insert phase and the lookup phase. Since we set the parameter values at the beginning of the workload execution and do not change it dynamically during runtime, the tiering system does not perform optimally for both the phases. Future tiering systems should identify and adapt to the program phases to achieve maximum performance.

**Intelligent tiering systems.** In our work, the optimizer was able to learn about the workload behavior and the appropriate parameter values by running multiple iterations of the workload. This could be replaced by a machine learning model that is trained on various workloads. In that case, the

model could be trained to predict the right parameter values or even the exact placement and migration decisions for a given workload.

## 6 Related Work

In this section, we discuss some of the research in memory tiering relevant to our work.

**Page-based memory tiering systems:** Most of the previous proposed tiering systems perform page migrations transparently to the applications. These systems are based on heuristics and use static thresholds to make policy decisions. Some use recency-based policies [28, 46] while others use frequency-based policies [3, 13, 23, 35, 47]. Some other use both page access recency and frequency information to make data placement decisions [13, 17, 27]. We find that these systems are sub-optimal as they do not properly adapt to the workload or the underlying hardware.

**Dynamic tuning of knobs in tiering systems:** There has been some prior work that try to adapt parts of a tiering system to the workload. Memtis [24] uses dynamic threshold adaptation for page hotness classification which leads to better fast memory tier utilization. Cori [10] tunes the periodicity of data movement in hybrid memory systems by extracting data reuse patterns from the application. Researches have proposed systems that migrate pages at different granularities for different workloads [2, 36]. Our research finds that it is important to adapt all parts of a tiering system for maximum performance. Instead of adjusting the parameters, some tiering systems use different policies for different workloads. Heo et.al., propose a dynamic policy selection mechanism which identifies the best migration policy amongst LRU, LFU and random for a given workload [17]. Yu et.al., build bandwidth-aware tiering systems [48].

**Machine learning for data placement.** Sibyl [38] uses reinforcement learning for data placement in hybrid storage systems. It observes different features of the running workload and storage devices such as number of accesses, access interval and free space to make data placement decisions. Kleio [9] uses deep neural networks to make intelligent page placement decision.

**Hardware tiering and profile-guided data placement.** To overcome the inaccuracies and inefficiencies in software profiling, Ramos et.al., propose augmenting the memory controller hardware to monitor access pattern and migrate pages between memories which would be more efficient [34].

X-Mem [11] and Mira [16] leverage profile guided techniques to determine object hotness offline and make data placement decisions during allocation time. These approaches do not work well for all applications, especially for those whose hot set changes over time.

**Optimizing knobs for systems.** Determining the right values for knobs has been a long-standing research problem in many systems. For example, LlamaTune [22] uses domain knowledge to efficiently sample and tune the knobs

for DBMS with Bayesian Optimization. Cao et al. [6] study the black-box auto-tuning for storage systems and find that optimal configurations are highly dependent on hardware, software, and workloads; no single technique is superior to all others. Our work is the first to shed light on optimizing knobs for memory-tiering systems.

## 7 Conclusion

Existing tiering solutions that use heuristics and static thresholds do not perform well under all scenarios. The key insight from our work is that parameters (knobs) of existing tiering systems can be tweaked to express different data placement and migration patterns. Tuning parameter values can yield significant performance benefits for most workloads running on different types of hardware. We find that a Bayesian optimizer can tune the parameters effectively as it able to successfully generate a mental model of the program behavior. We observe that in almost all cases, the optimizer identifies the workload pattern correctly. By using an optimizer to tune the parameter knobs of existing tiering systems, we can achieve up to 2x performance improvement.

## References

[1] Reto Achermann and Ashish Panwar. 2020. Mitosis Workload Btree. https://github.com/mitosis-project/mitosis-workload-btree.

[2] Shashank Adavally and Krishna Kavi. 2021. Subpage migration in heterogeneous memory systems. In *Workshop on Heterogeneous Memory Systems (HMEM-2021), Colocated with ICS 2021*.

[3] Neha Agarwal and Thomas F Wenisch. 2017. Thermostat: Application-transparent page management for two-tiered main memory. In *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems*. 631–644.

[4] Scott Beamer, Krste Asanović, and David Patterson. 2015. The GAP benchmark suite. *arXiv preprint arXiv:1508.03619* (2015).

[5] Zhen Cao, Geoff Kuenning, and Erez Zadok. 2020. Carver: Finding Important Parameters for Storage System Tuning. In *18th USENIX Conference on File and Storage Technologies (FAST 20)*. USENIX Association, Santa Clara, CA, 43–57. https://www.usenix.org/conference/fast20/presentation/cao-zhen

[6] Zhen Cao, Vasily Tarasov, Sachin Tiwari, and Erez Zadok. 2018. Towards Better Understanding of Black-box {Auto-Tuning}: A Comparative Analysis for Storage Systems. In *2018 USENIX Annual Technical Conference (USENIX ATC 18)*. 893–907.

[7] Intel Corporation. 2018. Intel 64 and IA-32 Architectures Software Developer Manuals. https://software.intel.com/articles/intel-sdm. (2018).

[8] Ian Cutress and Billy Tallis. 2018. Intel Launches Optane DIMMs up to 512GB: Apache Pass is Here. https://www.anandtech.com/show/12828/intel-launches-optane-dimms-up-to-512gb-apache-pass-is-here.

[9] Thaleia Dimitra Doudali, Sergey Blagodurov, Abhinav Vishnu, Sudhanva Gurumurthi, and Ada Gavrilovska. 2019. Kleio: A hybrid memory page scheduler with machine intelligence. In *Proceedings of the 28th International Symposium on High-Performance Parallel and Distributed Computing*. 37–48.

[10] Thaleia Dimitra Doudali, Daniel Zahka, and Ada Gavrilovska. 2021. Cori: Dancing to the right beat of periodic data movements over hybrid memory systems. In *2021 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. IEEE, 350–359.

[11] Subramanya R Dulloor, Amitabha Roy, Zheguang Zhao, Narayanan Sundaram, Nadathur Satish, Rajesh Sankaran, Jeff Jackson, and Karsten Schwan. 2016. Data tiering in heterogeneous memory systems. In *Proceedings of the Eleventh European Conference on Computer Systems*. 1–16.

[12] Dmitry Duplyakin, Robert Ricci, Aleksander Maricq, Gary Wong, Jonathon Duerig, Eric Eide, Leigh Stoller, Mike Hibler, David Johnson, Kirk Webb, Aditya Akella, Kuangching Wang, Glenn Ricart, Larry Landweber, Chip Elliott, Michael Zink, Emmanuel Cecchet, Snigdhaswin Kar, and Prabodh Mishra. 2019. The Design and Operation of CloudLab. In *2019 USENIX Annual Technical Conference (USENIX ATC 19)*. USENIX Association, Renton, WA, 1–14. https://www.usenix.org/conference/atc19/presentation/duplyakin

[13] Padmapriya Duraisamy, Wei Xu, Scott Hare, Ravi Rajwar, David Culler, Zhiyi Xu, Jianing Fan, Christopher Kennelly, Bill McCloskey, Danijela Mijailovic, et al. 2023. Towards an adaptable systems architecture for memory tiering at warehouse-scale. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*. 727–741.

[14] Ayat Fekry, Lucian Carata, Thomas Pasquier, Andrew Rice, and Andy Hopper. 2020. To Tune or Not to Tune? In Search of Optimal Configurations for Data Analytics. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Virtual Event, CA, USA) *(KDD '20)*. Association for Computing Machinery, New York, NY, USA, 2494–2504. https://doi.org/10.1145/3394486.3403299

[15] Johannes Freischuetz, Konstantinos Kanellis, Brian Kroth, and Shivaram Venkataraman. 2025. TUNA: Tuning Unstable and Noisy Cloud Applications. In *Proceedings of the Twentieth European Conference on Computer Systems* (Rotterdam, Netherlands) *(EuroSys '25)*. Association for Computing Machinery, New York, NY, USA, 954–973. https://doi.org/10.1145/3689031.3717480

[16] Zhiyuan Guo, Zijian He, and Yiying Zhang. 2023. Mira: A program-behavior-guided far memory system. In *Proceedings of the 29th Symposium on Operating Systems Principles*. 692–708.

[17] Taekyung Heo, Yang Wang, Wei Cui, Jaehyuk Huh, and Lintao Zhang. 2020. Adaptive page migration policy with huge pages in tiered memory systems. *IEEE Trans. Comput.* 71, 1 (2020), 53–68.

[18] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. 2011. Sequential model-based optimization for general algorithm configuration. In *Learning and Intelligent Optimization: 5th International Conference, LION 5, Rome, Italy, January 17-21, 2011. Selected Papers 5*. Springer, 507–523.

[19] Compute Express Link Consortium. Inc. 2020. CXL® Specification. https://computeexpresslink.org/cxl-specification/.

[20] Joseph Izraelevitz, Jian Yang, Lu Zhang, Juno Kim, Xiao Liu, Amirsaman Memaripour, Yun Joon Soh, Zixuan Wang, Yi Xu, Subramanya R Dulloor, et al. 2019. Basic performance measurements of the intel optane DC persistent memory module. *arXiv preprint arXiv:1903.05714* (2019).

[21] Konstantinos Kanellis, Ramnatthan Alagappan, and Shivaram Venkataraman. 2020. Too many knobs to tune? towards faster database tuning by pre-selecting important knobs. In *12th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 20)*.

[22] Konstantinos Kanellis, Cong Ding, Brian Kroth, Andreas Müller, Carlo Curino, and Shivaram Venkataraman. 2022. LlamaTune: Sample-efficient DBMS configuration tuning. *arXiv preprint arXiv:2203.05128* (2022).

[23] KyungSoo Lee, Sohyun Kim, Joohee Lee, Donguk Moon, Rakie Kim, Honggyu Kim, Hyeongtak Ji, Yunjeong Mun, and Youngpyo Joo. 2024. Improving key-value cache performance with heterogeneous memory tiering: A case study of CXL-based memory expansion. *IEEE Micro* (2024).

[24] Taehyung Lee, Sumit Kumar Monga, Changwoo Min, and Young Ik Eom. 2023. MEMTIS: Efficient Memory Tiering with Dynamic Page Classification and Page Size Determination. In *Proceedings of the 29th Symposium on Operating Systems Principles*. 17–34.

[25] Baptiste Lepers and Willy Zwaenepoel. 2023. Johnny Cache: the End of {DRAM} Cache Conflicts (in Tiered Main Memory Systems). In *17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23)*. 519–534.

[26] Mark Mansi, Bijan Tabatabai, and Michael M Swift. 2022. {CBMM}: Financial Advice for Kernel Memory Managers. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*. 593–608.

[27] Adnan Maruf, Ashikee Ghosh, Janki Bhimani, Daniel Campello, Andy Rudoff, and Raju Rangaswami. 2022. MULTI-CLOCK: Dynamic Tiering for Hybrid Memory Systems.. In *HPCA*. 925–937.

[28] Hasan Al Maruf, Hao Wang, Abhishek Dhanotia, Johannes Weiner, Niket Agarwal, Pallab Bhattacharya, Chris Petersen, Mosharaf Chowdhury, Shobhit Kanaujia, and Prakash Chauhan. 2023. TPP: Transparent page placement for CXL-enabled tiered-memory. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*. 742–755.

[29] Dingheng Mo, Fanchao Chen, Siqiang Luo, and Caihua Shan. 2023. Learning to Optimize LSM-trees: Towards A Reinforcement Learning based Key-Value Store for Dynamic Workloads. *Proc. ACM Manag. Data* 1, 3, Article 213 (Nov. 2023), 25 pages. https://doi.org/10.1145/3617333

[30] Timothy Morgan. 2020. CXL And Gen-Z Iron Out A Coherent Interconnect Strategy. https://www.nextplatform.com/2020/04/03/cxl-and-gen-z-iron-out-a-coherent-interconnect-strategy/. In *The Next Platform*.

[31] Richard C Murphy, Kyle B Wheeler, Brian W Barrett, and James A Ang. 2010. Introducing the graph 500. *Cray Users Group (CUG)* 19, 45-74 (2010), 22.

[32] SeongJae Park, Yunjae Lee, and Heon Y Yeom. 2019. Profiling dynamic data access patterns with controlled overhead and quality. In *Proceedings of the 20th International Middleware Conference Industrial Track*. 1–7.

[33] Steven J Plimpton, Ron Brightwell, Courtenay Vaughan, Keith Underwood, and Mike Davis. 2006. A simple synchronous distributed-memory algorithm for the HPCC RandomAccess benchmark. In *2006 IEEE International Conference on Cluster Computing*. IEEE, 1–7.

[34] Luiz E Ramos, Eugene Gorbatov, and Ricardo Bianchini. 2011. Page placement in hybrid memory systems. In *Proceedings of the international conference on Supercomputing*. 85–95.

[35] Amanda Raybuck, Tim Stamler, Wei Zhang, Mattan Erez, and Simon Peter. 2021. Hemem: Scalable tiered memory management for big data applications and real nvm. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*. 392–407.

[36] Jee Ho Ryoo, Lizy K John, and Arkaprava Basu. 2018. A case for granularity aware page migration. In *Proceedings of the 2018 International Conference on Supercomputing*. 352–362.

[37] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. 2015. Taking the human out of the loop: A review of Bayesian optimization. *Proc. IEEE* 104, 1 (2015), 148–175.

[38] Gagandeep Singh, Rakesh Nadig, Jisung Park, Rahul Bera, Nastaran Hajinazar, David Novo, Juan Gómez-Luna, Sander Stuijk, Henk Corporaal, and Onur Mutlu. 2022. Sibyl: Adaptive and extensible data placement in hybrid storage systems using online reinforcement learning. In *Proceedings of the 49th Annual International Symposium on Computer Architecture*. 320–336.

[39] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. 2012. Practical bayesian optimization of machine learning algorithms. *Advances in Neural Information Processing Systems* 25 (2012).

[40] Yan Sun, Yifan Yuan, Zeduo Yu, Reese Kuper, Chihun Song, Jinghan Huang, Houxiang Ji, Siddharth Agarwal, Jiaqi Lou, Ipoom Jeong, et al. 2023. Demystifying cxl memory with genuine cxl-ready systems and devices. In *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*. 105–121.

[41] John R Tramm, Andrew R Siegel, Tanzima Islam, and Martin Schulz. 2014. XSBench-the development and verification of a performance abstraction for Monte Carlo reactor analysis. *The Role of Reactor Physics toward a Sustainable Future (PHYSOR)* (2014).

[42] Stephen Tu, Wenting Zheng, Eddie Kohler, Barbara Liskov, and Samuel Madden. 2013. Speedy transactions in multicore in-memory databases. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*. 18–32.

[43] Johannes Weiner, Niket Agarwal, Dan Schatzberg, Leon Yang, Hao Wang, Blaise Sanouillet, Bikash Sharma, Tejun Heo, Mayank Jain, Chunqiang Tang, et al. 2022. TMO: Transparent memory offloading in datacenters. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*. 609–621.

[44] Lingfeng Xiang, Zhen Lin, Weishu Deng, Hui Lu, Jia Rao, Yifan Yuan, and Ren Wang. 2024. MATRYOSHKA: Non-Exclusive Memory Tiering via Transactional Page Migration. *arXiv preprint arXiv:2401.13154* (2024).

[45] Qiumin Xu, Huzefa Siyamwala, Mrinmoy Ghosh, Tameesh Suri, Manu Awasthi, Zvika Guz, Anahita Shayesteh, and Vijay Balakrishnan. 2015. Performance analysis of NVMe SSDs and their implication on real world databases. In *Proceedings of the 8th ACM International Systems and Storage Conference*. 1–11.

[46] Zi Yan, Daniel Lustig, David Nellans, and Abhishek Bhattacharjee. 2019. Nimble page management for tiered memory systems. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 331–345.

[47] Zhengyu Yang, Morteza Hoseinzadeh, Allen Andrews, Clay Mayers, David Thomas Evans, Rory Thomas Bolt, Janki Bhimani, Ningfang Mi, and Steven Swanson. 2017. AutoTiering: Automatic data placement manager in multi-tier all-flash datacenter. In *2017 IEEE 36th International Performance Computing and Communications Conference (IPCCC)*. IEEE, 1–8.

[48] Seongdae Yu, Seongbeom Park, and Woongki Baek. 2017. Design and implementation of bandwidth-aware memory placement and migration policies for heterogeneous memory systems. In *Proceedings of the International Conference on Supercomputing*. 1–10.

[49] Lei Zhang, Reza Karimi, Irfan Ahmad, and Ymir Vigfusson. 2020. Optimal data placement for heterogeneous cache, memory, and storage systems. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* 4, 1 (2020), 1–27.

[50] Chenxingyu Zhao, Tapan Chugh, Jaehong Min, Ming Liu, and Arvind Krishnamurthy. 2022. Dremel: Adaptive Configuration Tuning of RocksDB KV-Store. *Proc. ACM Meas. Anal. Comput. Syst.* 6, 2, Article 37 (June 2022), 30 pages. https://doi.org/10.1145/3530903