

Unified Holistic Memory Management Supporting Multiple Big Data Processing Frameworks over Hybrid Memories

LEI CHEN and JIACHENG ZHAO, SKL Computer Architecture, ICT, CAS, China and University of Chinese Academy of Sciences, China

CHENXI WANG, University of California, Los Angeles, California

TING CAO, Microsoft Research, China

JOHN ZIGMAN, The University of Sydney, Australia

HARIS VOLOS, University of Cyprus, Cyprus

ONUR MUTLU, ETH Zürich, Switzerland

FANG LV, SKL Computer Architecture, ICT, CAS, China

XIAOBING FENG, SKL Computer Architecture, ICT, CAS, China and University of Chinese Academy of Sciences, China

GUOQING HARRY XU, University of California, Los Angeles, California

HUIMIN CUI, SKL Computer Architecture, ICT, CAS, China and University of Chinese Academy of Sciences, China

To process real-world datasets, modern data-parallel systems often require extremely large amounts of memory, which are both costly and energy inefficient. Emerging **non-volatile memory (NVM)** technologies offer high capacity compared to DRAM and low energy compared to SSDs. Hence, NVMs have the potential to fundamentally change the dichotomy between DRAM and durable storage in Big Data processing. However, most Big Data applications are written in *managed languages* and executed on top of a *managed*

This submission is based on authors' previous publication: *Chenxi Wang, Huimin Cui, Ting Cao, John Zigman, Haris Volos, Onur Mutlu, Fang Lv, Xiaobing Feng, and Guoqing Harry Xu. 2019. Panthera: Holistic memory management for big data processing over hybrid memories. In Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2019). ACM, New York, 347–362. DOI:<https://doi.org/10.1145/3314221.3314650>.*

In this article, we develop a dynamic monitoring technique which monitors data access patterns at runtime to help data replacement when the patterns can not be inferred statically. We apply this technique to QuickCached and show its effectiveness. We also analyze the situations where coarse-grained patterns are not enough to achieve good performance. We further leverage the profiling technique to do fine-grained data replacement. The experiment results show that our methods can achieve more energy reduction with less performance overhead.

The work is supported in part by National Natural Science Foundation of China grants 62090024, 61872043, 61802368, and by US National Science Foundation grants CNS-1763172, CNS-1907352, CNS-2006437, CNS-2007737, CNS-2128653, and CNS-2106838, ONR grants N00014-16-1-2913 and N00014-18-1-2037, as well as gifts from Alibaba, Intel, and VMware.

Authors' addresses: L. Chen, J. Zhao (corresponding author), X. Feng, and H. Cui, SKL Computer Architecture, ICT, CAS, Beijing, China and University of Chinese Academy of Sciences, Beijing, China; email: huimin.cui@gmail.com; C. Wang and G. H. Xu, University of California, Los Angeles, USA; T. Cao, Microsoft Research, China; J. Zigman, The University of Sydney, Australia; H. Volos, University of Cyprus, Cyprus; O. Mutlu, ETH Zürich, Switzerland; F. Lv, SKL Computer Architecture, ICT, CAS, China.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

0734-2071/2022/07-ART2 \$15.00

<https://doi.org/10.1145/3511211>

runtime that already performs various dimensions of memory management. Supporting hybrid physical memories adds a new dimension, creating unique challenges in data replacement. This article proposes Panthera, a *semantics-aware, fully automated* memory management technique for Big Data processing over hybrid memories. Panthera analyzes user programs on a Big Data system to infer their coarse-grained access patterns, which are then passed to the Panthera runtime for efficient data placement and migration. For Big Data applications, the coarse-grained data division information is accurate enough to guide the GC for data layout, which hardly incurs overhead in data monitoring and moving. We implemented Panthera in OpenJDK and Apache Spark. Based on Big Data applications' memory access pattern, we also implemented a new profiling-guided optimization strategy, which is *transparent* to applications. With this optimization, our extensive evaluation demonstrates that Panthera reduces energy by 32–53% at less than 1% time overhead on average. To show Panthera's applicability, we extend it to QuickCached, a pure Java implementation of Memcached. Our evaluation results show that Panthera reduces energy by 28.7% at 5.2% time overhead on average.

CCS Concepts: • **Information systems** → **Data management systems**; • **Hardware** → **Non-volatile memory**; • **Software and its engineering** → **Memory management**;

Additional Key Words and Phrases: Hybrid memories, Big Data systems, memory management, garbage collection

ACM Reference format:

Lei Chen, Jiacheng Zhao, Chenxi Wang, Ting Cao, John Zigman, Haris Volos, Onur Mutlu, Fang Lv, Xiaobing Feng, Guoqing Harry Xu, and Huimin Cui. 2022. Unified Holistic Memory Management Supporting Multiple Big Data Processing Frameworks over Hybrid Memories. *ACM Trans. Comput. Syst.* 39, 1–4, Article 2 (July 2022), 38 pages.

<https://doi.org/10.1145/3511211>

1 INTRODUCTION

Modern Big Data computing exemplified by systems such as Spark and QuickCached is extremely memory intensive. Lack of memory can lead to a range of severe functional and performance issues including out-of-memory crashes, significantly degraded efficiency, or even loss of data upon node failures. Relying completely on DRAM to satisfy the memory need of a data center is costly in many different ways—e.g., large-volume DRAM is expensive and energy inefficient; furthermore, DRAM's relatively small capacity dictates that a large number of machines is often needed just to provide sufficient memory, resulting in underutilized CPU resources for workloads that cannot be easily parallelized.

Emerging non-volatile memory (NVM), such as phase change memory (PCM) [49, 79, 89], resistive random-access memory (RRAM) [78], Spin-transfer torque memory (STT-MRAM) [46] or 3D XPoint [5], is a promising technology that, has large memory capacity, energy efficiency and low per-GB cost, making them a supplement to traditional DRAM. NVM is on the memory bus and can be accessed via load/store instructions, enabling direct manipulation of persistent data in memory. There are two representative usages of NVM. The first is to leverage its persistence feature to ensure the persistent data structures are crash consistent and resume executions in the event of a failure [14, 21, 23, 24, 26, 27, 41, 42, 47, 48, 51, 52, 55, 56, 64, 71, 73, 81–84, 95, 96]. The second is as a supplement to traditional DRAM devices, i.e., to build the hybrid memory architecture, which benefits from both access speed of DRAM and the large capacity, low power consumption, and low per-GB cost of NVM. Our proposed approach falls into the second category. Systems with hybrid memories have received much attention [9, 11, 13, 15, 18, 25, 34, 43, 48, 49, 54, 60, 62, 63, 69–71, 76, 77, 80, 85–87, 90, 93, 94] recently from both academia and industry. The benefit of mixing NVM with DRAM for Big Data systems is obvious—NVM's high capacity makes it possible to fulfill the

high memory requirement of a Big Data workload with a small number of compute nodes, holding the promise of significantly reducing the costs of both hardware and energy in large data centers.

1.1 Problems

Although using NVM for Big Data systems is a promising direction, the idea has not yet been fully explored. Adding NVM naïvely would lead to large performance penalties due to its significantly increased access latency and reduced bandwidth—e.g. the latency of an NVM read is 2–4× larger than that of a DRAM read and NVM’s bandwidth is about 1/8–1/3 of that of DRAM [30, 75]. Hence, a critical research question that centers around all hybrid-memory-related research is *how to perform intelligent data allocation and migration between DRAM and NVM so that we can maximize the overall energy efficiency while minimizing the performance overhead?* To answer this question in the context of Big Data processing, there are two major challenges.

1.1.1 Challenge #1: Working with Garbage Collection (GC). A common approach to managing hybrid memories is to modify the OS or hardware to (1) monitor access frequency of physical memory pages and (2) move the hot (frequently accessed) data into DRAM. This approach works well for native language applications where data stays in the memory location it is allocated into. However, in *managed languages*, the garbage collector keeps changing the data layout in memory by copying objects to different physical memory pages, which breaks the bonding between data and physical memory address. Most Big Data systems are written in such managed languages, e.g., Java and Scala, for the quick development cycle and rich community support they provide. Managed languages are executed on top of a managed runtime such as the JVM, which employs a set of sophisticated memory management techniques such as garbage collection. As a traditional garbage collector is not aware of hybrid memories, allocating and migrating hot/cold pages at the OS level can easily lead to interference between these two different levels of memory management.

1.1.2 Challenge #2: Working with Application-Level Memory Subsystems. Modern Big Data systems all contain sophisticated memory subsystems that perform various memory management tasks *at the application level*. For instance, Apache Spark [6] uses **resilient distributed datasets (RDDs)** as its data abstraction. An RDD is a distributed data structure that is partitioned across different servers. At a low level, each RDD partition is an array of Java objects, each representing a data tuple. RDDs are often immutable but can exhibit diverse lifetime behavior. For example, developers can explicitly persist RDDs in memory for memorization or fault tolerance. Such RDDs are long-lived while RDDs storing intermediate results are short-lived.

An RDD can be at one of many storage levels (e.g., memory, disk, unmaterialized, etc.). Spark further allows developers to specify, with annotations, where an RDD should be allocated, e.g., in the managed heap or native memory. Objects allocated natively are not subject to GC, leading to increased efficiency. However, data processing tasks, such as shuffle, join, map, or reduce, are performed over the managed heap. A native-memory-based RDD cannot be directly processed unless it is first moved into the heap. Hence, where to allocate an RDD depends on when and how it is processed. For example, a frequently accessed RDD should be placed in DRAM while a native-memory-based RDD would not be frequently used and placing it in NVM would be desirable. Clearly, efficiently using hybrid memories requires appropriate coordination between these orthogonal data placement policies, i.e., the heap, native memory, or disk, vs. NVM or DRAM.

Another instance is QuickCached, which is a pure Java implementation of Memcached server based on QuickServer [3]. In particular, it is an in-memory key-value store for small chunks of arbitrary data (strings, objects) from results of database calls, API calls, or page rendering. QuickCached leverages *ConcurrentHashMap* and *SoftReference* for managing its memory, i.e.,

ConcurrentHashMap for storing the key-value data, and *SoftReference* for automatically clearing data at the discretion of the garbage collector in response to memory demand.

1.1.3 State of the Art. In summary, the key challenges in supporting hybrid memories for Big Data processing lie in how to develop runtime system techniques that can make memory allocation/migration decisions that match how data is actually used in an application. Although techniques such as Espresso [80] and Write Rationing [11] support NVM for managed programs, neither of them was designed for Big Data processing whose data usage is greatly different than that of regular, non-data-intensive Java applications [65, 66].

For example, Espresso defines a new programming model that can be used by the developer to allocate objects in persistent memory. However, real-world developers would be reluctant to completely re-implement their systems from scratch using such a new model. Shoaib et al. [11] introduced the Write Rationing GC, which moves the objects that experience a large/small number of writes into DRAM/NVM to prolong NVM's lifetime. Write Rationing pioneers the work of using the GC to migrate objects based on their access patterns. However, Big Data systems make heavy use of immutable datasets—for example, in Spark, most RDDs are immutable. Placing all immutable RDDs into NVM can incur a large overhead as many of these RDDs are frequently read and an NVM read is 2–4× slower than a DRAM read.

1.2 Our Contributions

1.2.1 Our Insight. We analyzed two representative big data systems, i.e., Spark for data processing and QuickCached for data store, and we observed that even they exhibit diverse memory behaviors, we have opportunities to share the common memory management strategy in JVM. Our observations are as follows:

- Spark applications have two unique characteristics that can greatly aid hybrid memory management.

First, they perform bulk object creation, and data objects exhibit strong *epochal behavior and clear access patterns*. For example, Spark developers program with RDDs, each of which contains objects with exactly the same access/lifetime patterns. Exploiting these patterns at the runtime would make it much easier for Big Data applications to enjoy the benefits of hybrid memory.

Second, the data access and lifetime patterns are often *statically* observable in the *user program*. For example, an RDD is a coarse-grained data abstraction in Spark and the access patterns of different RDDs can often be inferred from the way they are created and used in the program (Section 2).

- QuickCached has one unique characteristic that can aid hybrid memory management. It uses a huge hash table for storing the key-value pairs, and when processing each query request, it would create a large number temporary objects which would be frequently accessed in a very short period of time. Therefore, the lifetime of frequently accessed data is short, and the data with long lifetime are infrequently accessed.

Hence, unlike regular, non-data-intensive applications for which profiling is often needed to understand the access patterns of individual objects, we can develop a simple static analysis for a Big Data application to infer the access pattern of each coarse-grained data collection, in which all objects share the same pattern. This observation aligns well with prior work (e.g., Facade [66] or Yak [65]) that requires simple annotations to specify epochs to perform efficient garbage collection for Big Data systems. The static analysis does not incur any runtime overhead, yet it can produce precise enough data access information for the runtime system to perform effective allocation and migration.

1.2.2 Panthera. Based on our extensive experience with Big Data applications, we propose Panthera, which divides a mess of data objects into several data collections according to application's semantics and infers the coarse-grained data usage behavior by lightweight static program analysis and dynamic data usage monitoring. Panthera leverages garbage collection to migrate data between DRAM and NVM, incurring almost no runtime overhead.

We select two big data processing frameworks in this article. First, we focus on Apache Spark as it is the de-facto data-parallel framework deployed widely in industry. Spark hosts a range of applications in machine learning, graph analytics, stream processing, and so on, making it worthwhile to build a specialized runtime system, which can provide immediate benefit to all applications running atop. Furthermore, to demonstrate the generality of our approach, Panthera is built also on QuickCached, a Java implementation of Memcached, and Section 4 provides a detailed discussion of Panthera's applicability.

Panthera enhances both the JVM and Spark/QuickCached with two major innovations. First, based on the observation that access patterns in a Big Data application can be identified statically, we develop two static analyzers (Section 3) for Spark and QuickCached, respectively. In particular, the Spark analyzer analyzes a Spark program to infer a memory tag (i.e., NVM or DRAM) for each RDD variable based on the variable's location and the way it is used in the program, and the QuickCached analyzer analyzes the QuickCached source code to identify the huge global hash table, and then infer its corresponding memory tags. These tags indicate which memory the objects should be allocated in.

Second, we develop a new semantics-aware and physical-memory-aware generational GC (Section 4). Our static analysis instruments the Spark program and QuickCached to pass the inferred memory tags down to the runtime system, which uses these tags to make allocation/migration decisions. Since our GC is based on a high-performance generational GC in OpenJDK, Panthera's heap has two spaces, representing a young and an old generation. We place the entire young generation in DRAM while splitting the old generation into a small DRAM component and a large NVM component. The insight driving this design is based on a set of key observations (discussed in Section 2 in detail) we make over the lifetimes and access patterns of RDDs in representative Spark executions and the QuickCached objects:

- Most objects are allocated initially in the young generation. Since they are frequently accessed during initialization, placing them in DRAM enables fast access to them.
- Long-lived objects in Spark can be roughly classified into two categories: (1) long-lived RDDs that are frequently accessed during data transformation (e.g., cached for iterative algorithms) and (2) long-lived RDDs that are cached primarily for fault tolerance. The first category of RDDs should be placed in the DRAM component of the old generation because they have long lifespans and DRAM provides desirable performance for frequent access to them. The second category should be placed in the NVM component of the old generation because they are infrequently accessed and hence NVM's large access latency has relatively small impact on overall performance.
- For Spark programs, there are also short-lived RDDs that store temporary, intermediate results. These RDDs die and are then reclaimed in the young generation quickly, leading to frequent accesses to this area. This is another reason why we place the young generation within DRAM.
- For QuickCached, there is only one long-lived object, i.e., *ConcurrentHashMap* for storing the key-value data. Among the hash table, only a small fraction is frequently accessed for a specific request, which will be identified at runtime with negligible overhead. Thus, the *ConcurrentHashMap* should be placed in the NVM component of the old generation, except the identified frequently accessed fraction.

```

Top: obj org/apache/spark/rdd/ShuffledRDD
depth [0]: array , [Lscala/Tuple2;
depth [1]: obj scala/Tuple2
depth [2]: obj java/lang/String
depth [3]: array [C
depth [2]: obj spark/util/collection/CompactBuffer
depth [3]: array , [Ljava/lang/String;
depth [4]: obj java/lang/String
depth [5]: array [C
depth [4]: obj java/lang/String
depth [5]: array [C

```

Fig. 1. The heap structure of an example RDD.

- For QuickCached, a number of temporary objects would be created to process a specific request. These objects are allocated in the young generation, and should be placed in DRAM enabling fast accesses.

Based on these observations, we modified both the minor and major GC, which allocate and migrate data objects, based on their RDD types and the semantic information inferred by our static analysis, into the spaces that best fit their lifetimes and access patterns. Our runtime system also monitors the transformations invoked over RDD objects to perform runtime (re)assessment of RDDs’ access patterns. Even if the static analysis does not accurately predict an RDD’s access pattern and the RDD gets allocated in an undesirable space, Panthera can still migrate the RDD from one space to another using the major GC.

1.2.3 Results. We have evaluated Panthera extensively with Spark applications, including graph computing (GraphX), machine learning (MLlib) and other iterative in-memory computing applications (Table 4), and QuickCached using Yahoo! Cloud Serving Benchmark (YCSB) [22]. Results with various heap sizes and DRAM ratios demonstrate that Panthera makes effective use of hybrid memories—overall, the Panthera-enhanced JVM reduces the memory energy by 22%–34% with only a 1%–9% execution time overhead for QuickCached, and reduces the memory energy by 32%–53% with only less than 1% execution time overhead on average for Spark, whereas Write Rationing [11] that moves read-only RDD objects into NVM incurs a 41% time overhead.

2 BACKGROUND AND MOTIVATION

This section provides necessary background on Apache Spark [6] and QuickCached [3] with motivating examples that illustrate the access patterns in a Spark program.

2.1 Spark Basics

Spark is a data-parallel system that supports acyclic data flow and in-memory computing. The major data representation used by Spark is **resilient distributed dataset (RDD)** [91], which represents a read-only collection of tuples. An RDD is a distributed memory abstraction partitioned in the cluster. Each partition is an array of data items of the same type. Each node maintains an RDD partition, which is essentially a multi-layer Java data structure—a top RDD object references a Java array, which, in turn, references a set of tuple objects such as key-value pairs. Figure 1 shows the heap structure for an example RDD where each element is a pair of a string (key) and a compact buffer (value).

A Spark pipeline consists of a sequence of *transformations* and *actions* over RDDs. A transformation produces a new RDD from a set of existing RDDs; examples are *map*, *reduce*, or *join*. An action is a function that computes statistics from an RDD, such as an aggregation. Spark leverages *lazy evaluation* for efficiency, that is, a transformation may not be evaluated until an action is

performed later on the resulting RDD. Before data processing starts, the dependences between RDDs are first extracted from the transformations to form a *lineage graph*, which can be used to conduct lazy evaluation and RDD recomputation upon node failures.

With lazy evaluation, a transformation only creates a (top-level) RDD object without *materializing* the RDD (i.e., the point at which its internal array and actual data tuples are created). Recomputing all RDDs is time-consuming when the lineage is long or when it branches out, and hence, Spark allows developers to cache certain RDDs in memory (by using the API `persist`). Developers can specify a storage level for a persisted RDD, e.g., in memory or on disk, in the serialized or deserialized form, and the like. RDDs that are *not* explicitly persisted are temporary RDDs that will be garbage-collected when they are no longer used, while persisted RDDs are materialized and never collected.

The Spark scheduler examines the lineage graph to build a DAG of stages for execution. The lineage (transformation)-based dependences are classified into “narrow” and “wide”. A narrow dependence exists from a parent to a child RDD if each partition of the parent is used by *at most one* partition of the child RDD. By contrast, a wide dependence exists when each partition of the parent RDD may be used by *multiple* child partitions. Distinguishing these two types of dependences makes it possible for Spark to determine whether a shuffle is necessary. For example, for narrow dependences shuffling is not necessary, while for wide dependences it is.

A Spark pipeline is split into a set of *stages* based on shuffles (and thus wide dependences)—each stage ends at a shuffle that writes RDDs onto the disk and the next stage starts by reading data from disk files. Transformations that exhibit narrow dependences are grouped into the same stage and executed in parallel.

2.2 RDD Characteristics

An RDD is, at a low level, an array of Java objects, which are managed by the semantics-agnostic GC in the JVM. RDDs often exhibit predictable lifetime and memory-access patterns. Our goal is to pass these patterns down to the GC, which can exploit such semantic information for efficient data placement. We provide a concrete example to illustrate these patterns.

Figure 2(a) shows the Spark program for PageRank [19], which is a well-known graph algorithm used widely by search engines to rank web pages. The program iteratively computes the rank of each vertex based on the contributions of its in-neighbors. Three RDDs can be seen from its source code: `links` representing edges from the input graph, `contribs` containing contributions from incoming edges of each vertex, and `ranks` that maps each vertex to its page rank. `links` is a static map computed from the input while `contribs` and `ranks` are recomputed per iteration of the loop.

In addition to these three developer-defined RDDs visible in the program, Spark generates many invisible RDDs to store intermediate results during execution. A special type of intermediate RDD is `ShuffledRDD`. Each iteration of the loop in the example forms a stage that ends at a shuffle, writing shuffled data into different disk files. In the beginning of the next stage, Spark creates a `ShuffledRDD` as input for the stage. Unlike other intermediate RDDs that are never materialized, `ShuffledRDD`s are immediately materialized because they contain data read freshly out of disk files. However, since they are not persisted, they will be collected when the stage is completed.

In summary, (1) persisted RDDs are materialized at the moment the method `persist` is called and (2) non-persisted RDDs are not materialized unless they are `ShuffledRDD`s or an action is invoked on them.

2.3 Example

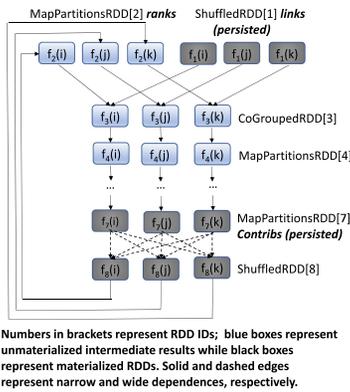
Figure 2(b) shows the set of RDDs that exists within a stage (i.e., iteration) and their dependences. Suppose each RDD has three partitions (on three nodes). The dashed edges represent wide

```

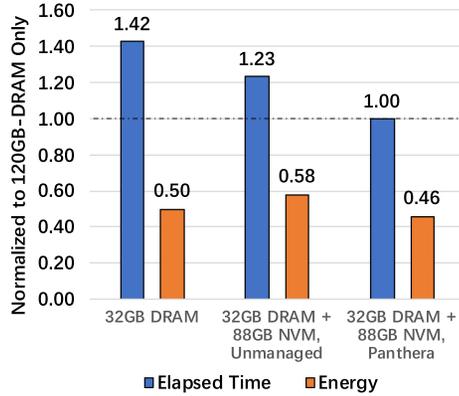
1 var lines = ctx.textFile(args[0], slices)
2 var links = lines.map{s=>
3   var parts = s.split("\\s+")
4   (parts(0), parts(1))
5 }.distinct().groupByKey()
6 .persist(StorageLevel.MEMORY_ONLY)
7
8 var ranks = links.mapValues(v => 1.0)
9 for(i <- 1 to iters){
10  var contribs = links.join(ranks).values.flatMap{
11    case(urls, rank) =>
12      val size = urls.size
13      urls.map(url=>(url, rank/size))
14        .persist(StorageLevel
15          .MEMORY_AND_DISK_SER)
16  }
17  ranks = contribs.reduceByKey(_ + _).
18    mapValues(0.15 + 0.85 * _)
19 }
20 ranks.count()

```

(a) PageRank program.



(b) Transformations within a stage.



(c) Results of DRAM-only and DRAM+NVM, managed by the OS and by Panthera.

Fig. 2. Characteristics of RDDs in Spark PageRank.

dependencies (i.e., shuffles) due to the reduction on Line 17. There are totally eight RDDs generated in each iteration. `ShuffledRDD[8]`, which stems from the reduction on Line 17, is transformed to `ranks` via a map transformation. `ranks` joins with `links` to form `CoGroupedRDD[3]`, which is then processed by four consecutive map functions, i.e., $f_4 - f_7$, producing `contribs` at the end. For unmaterialized (blue) RDDs, the sequence of transformations (e.g., $f_4 \circ \dots \circ f_7$) is applied to each record from the source RDD in a *streaming* manner via iterators to produce a final record.

For `links` and `contribs`, the developer invokes the method `persist` to materialize these RDDs. The storage levels indicate that `links` is cached in memory throughout the execution (as it is used in each iteration) while `contribs` generated in each iteration is kept in memory but will be serialized to disk upon memory pressure. `ranks` is not explicitly persisted. Hence, it is not materialized until the execution reaches Line 20 where action `count` is invoked on the RDD object.

The lifetime patterns of these different RDDs fall into two categories. Non-persisted intermediate RDDs are short-lived as their data objects are generated only during a pipelined execution.

Persisted RDDs are long-lived and stay in memory/on disk until the end of the execution. Their access patterns are, however, more diverse. Objects in an intermediate RDD are accessed at most once during streaming. Objects in a persisted RDD can exhibit different types of behavior. For RDDs like `links` that are used in each iteration, their objects are frequently accessed. In contrast, RDDs like `contribs` are persisted primarily for speeding up recovery from faults, and hence, their objects are rarely used after generated.

2.4 Design Choices

The different characteristics of DRAM and NVM make them suitable for different types of datasets. DRAM has low capacity and fast access speed, while NVM has large capacity but slow speed. Hence, DRAM is a good choice for storing small-sized, frequently accessed datasets, while large-sized, infrequently accessed datasets fit naturally into NVM. The clear distinction in the lifespans and access patterns of different RDDs makes it easy for them to be placed into different memories suitable for their behavior. For example, intermediate (blue) RDDs are never materialized. Their objects are created individually during streaming and then quickly collected by the GC. These objects are allocated in the young generation and will eventually die there. As a result, the memory used as young generation is frequently reused by these short-lived objects, which cause very high read/write frequency to this part of memory. This motivates our design choice of placing the young generation in DRAM, which matches the conclusion of previous works [11, 76].

Persisted RDDs, in contrast, have all their data objects created at the same time, and thus need large storage space. Since they are kept alive *indefinitely*, they should be allocated directly in the old generation. One category of persisted RDDs includes those that are frequently accessed, like `links`; they need to be placed in DRAM. Another category includes RDDs that are rarely accessed and cached for fault tolerance, like `contribs`, these RDDs should be placed in NVM. This behavioral difference motivates our choice of splitting the old generation into a DRAM and an NVM component.

We perform what we suggest on a system with 128-GB memory using Spark-based PageRank as the benchmark. For this experiment, we allocate 120-GB memory for the Spark and reserve 8-GB memory for the OS and other services. For the 120-GB Spark memory, 32-GB are DRAM and others are NVM. (We varied the DRAM ratios in evaluation section.) Figure 2(c) shows the performance and energy consumption normalized to a system with 120 GB of DRAM. Compared to using only 32-GB DRAM, adding 88-GB NVM to the system provides modest performance benefit (15%) but leads to 16% higher energy consumption, without proper data placement across DRAM and NVM (see *Unmanaged*, Section 7.2). After applying Panthera, RDD `links` and `contribs` are placed into DRAM and NVM, respectively. With such careful placement of data across DRAM and NVM, we find that (1) performance increases by 42% compared to using only a 32-GB DRAM, and becomes at the same level of the performance of using 120-GB DRAM; (2) energy consumption is 9% less than using only a 32-GB DRAM, and 54% less than using a 120-GB DRAM. We conclude that careful data placement between DRAM and NVM can provide the performance of large DRAM system, while keeping the energy consumption at the level of a small DRAM system.

2.5 QuickCached Basics

QuickCached is a pure Java implementation of Memcached server based on QuickServer, and it serves as an in-memory key-value store for small chunks of arbitrary data (strings, objects) from results of database calls, API calls, or page rendering.

QuickCached supports different backends for organizing and managing the key-value data, and its default backend leverages *ConcurrentHashMap* and *SoftReference*, i.e., *ConcurrentHashMap* for storing the key-value data, and *SoftReference* for automatically clearing data at the discretion

of the garbage collector in response to memory demand [3]. When processing each query request, QuickCached would create a large number of frequently accessed temporary objects, which will be destroyed when current query request is finished. Therefore, the lifetime of these frequently accessed data is short. In comparison, the *ConcurrentHashMap* provides full-lifecycle service and has long lifetime, however, only a small fraction in *ConcurrentHashMap* would be frequently accessed when processing one query request. Therefore, we have an opportunity to identify the data structure of *ConcurrentHashMap* statically, and identify its frequently accessed objects at runtime. Correspondingly, the *ConcurrentHashMap* should be placed in the NVM component of the old generation, and the identified frequently accessed objects should be placed in the DRAM component of the old generation.

3 STATIC INFERENCE OF MEMORY TAGS

Based on our observation of memory access patterns, we developed a simple static analysis that extracts necessary semantic information for efficient data placement. For Spark, the access patterns of RDDs can often be identified from the program using them. Our analysis automatically infers, for each persisted RDD visible in the program, whether it should be allocated in DRAM or NVM. This information is then passed down to the runtime system for appropriate data allocation. For QuickCached, our analysis takes user-annotated source codes as input and automatically generates code for runtime data placement, with the details discussed below.

3.1 Spark Analyzer

Static Analysis. In a Spark program, the developer can invoke `persist` with a particular storage level on an RDD to materialize the RDD, as illustrated in Figure 2. We piggyback on the storage levels to further determine if a persisted RDD should be placed into DRAM or NVM. In particular, Panthera statically analyzes the program to infer a memory tag (i.e., DRAM or NVM) for each `persist` call. Each of the ten existing storage levels (e.g., `MEMORY_ONLY`), except for `OFF_HEAP` and `DISK_ONLY`, is expanded into two sub-levels, annotated with NVM and DRAM, respectively (e.g., `MEMORY_ONLY_DRAM` and `MEMORY_ONLY_NVM`). `OFF_HEAP` is translated directly into `OFF_HEAP_NVM` because RDDs placed in native memory are rarely used, while `DISK_ONLY` does not carry any memory tag.

Our static analysis performs inference based on the *def-use* information w.r.t. each RDD variable declared in the program as well as the loop(s) in which the variable is defined/used. Our key insight is that if the variable is *defined* in each iteration of a computational loop, most of the RDD instances represented by the variable are *not* used frequently. This is because Spark RDDs are often immutable and hence, every definition of the RDD variable creates a new RDD instance at run time, leaving the old RDD instance cached and unused. Hence, we tag the variable “NVM”, instructing the runtime system to place these RDDs in NVM. An example is the `contribs` variable in Figure 2(a), which is defined in every iteration of the loop—although the variable is also used in each iteration, the use refers to the most recent RDD instance created in the last iteration while the instances created in all the other past iterations are left unused.

By contrast, if a variable is *used-only* (i.e., never defined) in the loop, such as `links`, we create a tag “DRAM” for it since only one instance of the RDD exists and is repeatedly used. Panthera analyzes not only RDD variables on which `persist` is explicitly called, but also those on which actions are invoked, such as the `ranks` variable in Figure 2(a). The tag inferred for an RDD variable (say v) is passed, at the materialization point of every RDD instance (v refers to), into the runtime system via automatically instrumented calls to auxiliary (native) methods provided by the Panthera JVM. We piggyback on a tracing GC to propagate this tag from the RDD object down to each data object contained in the RDD—when the GC runs, it moves objects with the same tag together into the same (DRAM or NVM) region (see Section 4).

One constraint that needs to be additionally considered is the location of the loop relative to the location of the materialization point of the RDD. We analyze the loop only if the materialization point *precedes or is in* the loop. Otherwise, whether the variable is used or defined in the loop does not matter as the RDD has not been materialized yet. For instance, although the `ranks` variable is defined in the loop that starts at Line 17, it does not get materialized until Line 20 after the loop finishes. Hence, its behavior in the loop does not affect its memory tag, which should actually depend on its *def-use* in the loops, if any, after Line 20.

If no loop exists in a program, the program has only one iteration and all RDDs receive an “NVM” tag as none of them are repeatedly accessed. If there are multiple loops to be considered for an RDD variable, we tag it “DRAM” if there exists one loop in which the variable is used-only and that loop follows or contains the materialization point of the RDD. The variable receives an “NVM” tag otherwise. If all persisted RDDs receive an “NVM” tag at the end of the analysis, we change the tags of all RDDs to “DRAM”—the goal is to fully utilize DRAM by first placing RDDs in DRAM. Once DRAM capacity is exhausted, the remaining RDDs, including those with a “DRAM” tag, will be placed in NVM.

Note that our analysis infers tags only for the RDD variables explicitly declared in the program. Intermediate RDDs produced during execution are not materialized and thus do not receive memory tags from our analysis. We discuss how to handle them in Section 4.

The memory tag of an RDD variable is a *static approximation* of its access pattern, which may not reflect the behavior of all RDD instances represented by the variable at run time. However, user code for data processing often has a simple batch-transformation logic. Hence, the static information inferred from our analysis is often good enough to help the runtime make an accurate placement decision for the RDD. In case the statically inferred tags do not precisely capture the RDD’s access information, Panthera has the ability to move RDDs between NVM and DRAM (within the old generation) based on their access frequencies, when a full-heap GC occurs. The dynamic data migration frequency is a good indicator for the accuracy of the static analysis. Section 4 provides a full discussion for this mechanism and Section 7.5 evaluate the accuracy of the static analysis and the overhead of dynamic migration.

Dealing with ShuffledRDD. Recall from Section 2 that, in addition to the RDDs on which `persist` is explicitly invoked, `ShuffledRDD`s, which are created from disk files after a shuffle, are also materialized. These RDDs are often the input of a stage but invisible in the program code. The challenge here is where to place them. Our insight is that their placement should depend on the other materialized RDDs that are transformed from (i.e., depend on) them in the same stage.

For example, in Figure 2(b), the input of the stage are two sets of `ShuffledRDD`s: [1] and [8]. `ShuffledRDD[1]` is the RDD represented by `links` and our static analysis already infers tag “DRAM” for it. `ShuffledRDD[8]` results from the reduction in the previous stage. Because `ShuffledRDD[8]` transitively produces `MapPartitionRDD[7]` (represented by `contribs`) and `MapPartitionRDD[7]` has a memory tag “NVM” inferred by our static analysis, we tag `ShuffledRDD[8]` “NVM” as well.

The main reason is that RDDs belonging to the same stage may share many data objects for optimization purposes. For example, a map transformation that only changes the values (of key-value pairs) in RDD *A* may generate a new RDD *B* that references the same set of key objects as in *A*. If *B* has already received a memory tag from our static analysis, it is better to assign the same tag to *A* so that these shared objects do not receive inconsistent tags and would not need to be moved from one memory to another when *B* is generated from *A*. This is especially beneficial when the transformation is in a computational loop—a large number of objects would be moved if *A* and *B* have different memory tags.

Figure 3 depicts our algorithm which assigns the same tag to *A* and *B*. We add support that scans the lineage graph at the beginning of each stage to propagate the memory tag *backward*,

```

1:  $G \leftarrow$  the lineage graph
2: while there are materialized RDDs not being selected do
3:    $V \leftarrow$  the lowest materialized RDD which has not been selected
4:   for all  $P \in \text{parent}(V)$  do
5:     if  $V.\text{tag} = \text{DRAM}$  then
6:        $P.\text{tag} \leftarrow \text{DRAM}$ 
7:     else if  $P.\text{tag} \neq \text{DRAM}$  then
8:        $P.\text{tag} \leftarrow V.\text{tag}$ 
9:     end if
10:  end for
11: end while
12: end

```

Fig. 3. Algorithm to assign same tags to RDDs which share data objects.

starting from the lowest materialized RDD in the graph that has received a tag from our analysis. Conflicts may occur during the propagation—an RDD encountered during the backward traversal may have an existing tag that is different from the tag being propagated. To resolve conflicts, we define the following priority order: DRAM > NVM, which means that upon a conflict, the resulting tag is always DRAM. This is because our goal is to minimize the NVM-induced overhead; RDDs with a “DRAM” tag inferred will be frequently used and putting them in NVM would cause large performance degradation.

3.2 QuickCached Analyzer

For QuickCached, the core object is the storage object, i.e., the hash table *ConcurrentHashMap*, and our QuickCached analyzer requires users annotate this object using the following syntax, @CoreHashObject, e.g.,

```

@CoreHashObject
ConcurrentHashMap hashTable;

```

According to the observation that only a small fraction of the hash table will be frequently accessed during the processing of one query request, thus the annotation will guide Panthera to place the hash table in NVM, meanwhile keep only the frequently accessed fraction in DRAM. In particular, the QuickCached analyzer identifies the annotated hash table *ConcurrentHashMap* and infer its memory tag as “NVM” statically.

However, different with Spark applications, the frequently accessed data in *ConcurrentHashMap* cannot be statically identified, since it is determined by the incoming requests at runtime. Thus, the static analysis is inefficient to infer meaningful tags for the objects stored in *ConcurrentHashMap*. To address this problem, the QuickCached analyzer introduced a dynamic mechanism to distinguish the frequently accessed data and tag it as “DRAM” at runtime. Since data accesses are highly skewed in real-world workloads [16], the frequently accessed data can be identified by monitoring the data access patterns at runtime. In particular, the analyzer automatically inserts some instrumentation codes at the callsites of the *get* method which is the interface for accessing the annotated *ConcurrentHashMap*. The instrumented codes are shown in Figure 4, serving to leverage a simple LRU strategy to tag the most recently accessed *value* data as “DRAM”.

In Spark, the RDD is an abstraction which is an array of Java objects at a low level, and such semantics facilitates the memory tags analysis and passing in Panthera. However, QuickCached lacks of such RDD abstraction, therefore, to share the same memory tag passing mechanism with Spark, we synthesize an RDD to wrap the objects that need to be managed by Panthera.

```

1 rdd_indicator("DRAM");
2 Object [] AggDramValues = new Object [MAX_OBJECTS_NUMBER];
3 long index = 0L;
4 ...
5 public Object get(...){
6     ...
7     /* retrieve value from hash table */
8     Object value = hashTable[key];
9     /* aggregate the most recently accessed data */
10    AggDramValues[(index++) % MAX_OBJECTS_NUMBER] = value;
11    ...
12 }

```

Fig. 4. Example of QuickCached analyzers instrumentation codes.

In particular, the instrumented codes work as follows. First, `rdd_indicator("DRAM")` ((Line 1)) declares the synthesized RDD in QuickCached, which will behave as the RDD in Spark. Second, we allocate a fixed-size auxiliary object array `AggDramValues` which will be wrapped in the synthesized RDD, to aggregate the most recently accessed `value` data together (Lines 2 and 3). Finally, the most recently visited elements would be copied into `AggDramValues` when they are accessed (Line 10).

With the synthesized RDD, Panthera provides a unified memory-tag-passing mechanism which can support both Spark and QuickCached, as will be discussed in Section 4.2.

4 THE PANTHERA GARBAGE COLLECTOR

While our static analysis (Section 3) determines where RDDs should be allocated, this information has to be communicated down to the runtime system, which recognizes only objects, not RDDs. Hence, our goal is to develop a new GC that, when placing/moving data objects, is aware of (1) the high-level semantics about where (DRAM or NVM) these RDDs should be placed and (2) the low-level information about the RDDs to which these objects belong.

We have implemented our new collection algorithm in OpenJDK 8 (build `jdk8u76-b02`) [8]. In particular, we have modified the object allocator, the interpreter, the two JIT compilers (C1 and Opto), and the Parallel Scavenge collector.

4.1 Design Overview

Heap Design. The Panthera GC is based on the Parallel Scavenge collector, which is the default GC in OpenJDK8. The collector divides the heap into a young and an old generation. As discussed earlier in Section 1, Panthera places the young generation in DRAM and splits the old generation into a DRAM component and an NVM component. The off-heap native memory is placed entirely in NVM. We reserve two unused bits, referred to as `MEMORY_BITS`, from the header of each object to indicate whether the object should be allocated into DRAM (01) or NVM (10). The default value for these bits is 00—objects that do not receive a tag have this default value. They will be promoted to the NVM component of the old generation if they live long enough. Figure 5 illustrates the heap structure and our allocation policies.

Allocation Policies. As discussed in Section 3, each materialized RDD carries a memory tag that comes from our static analysis or lineage-based tag propagation. However, at a low level, an RDD is a structure of objects, as illustrated in Figure 1, and these objects are created at different points in the execution. Our goal is to place all objects belonging to the same logical RDD—including the top object, the array object, tuple objects, and other objects reachable from tuples—together in

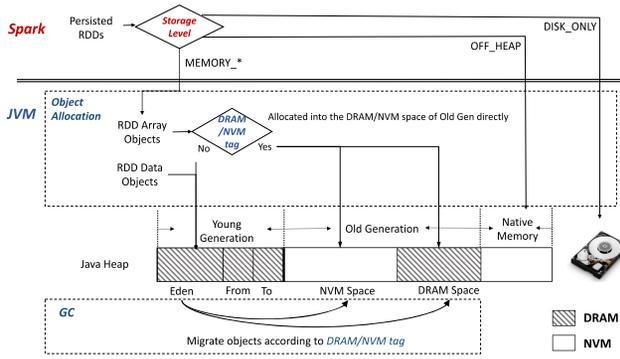


Fig. 5. The Panthera heap and allocation policies. Here RDD array objects refer to RDDs’ backbone arrays while data objects refer to other non-array objects in an RDD structure.

the space suggested by the RDD’s memory tag, because these objects likely have the same access pattern and lifetime.

However, this is rather challenging—our static analysis infers a memory tag for each *top RDD object* (whose type is a subtype of `org.apache.spark.rdd.RDD`) in the user program and we do not know what other objects belong to this RDD by just analyzing the user program. Statically identifying what objects belong to a logical data structure would require precise context-sensitive static analysis of *both user and system code*, which is difficult to do due to Spark’s extremely large codebase and the scalability issues of static analysis.

Our idea to solve this problem is that instead of attempting to allocate all objects of an RDD directly into the space (say *S*) suggested by the RDD’s tag, we *allocate only the array object into S upon its creation*. This is much easier to do—Panthera instruments each materialization point (e.g., before a call to `persist` or a Spark action) in the user program to pass the tag down to the runtime system without needing to analyze the Spark system code. Since the array is created at materialization, the runtime system can just use the tag to determine where to place it. All other objects in the RDD are not immediately allocated in *S* due to the aforementioned difficulties in finding their allocation sites. They are instead allocated in the young generation. Later, we use the GC to move these objects into *S* as tracing is performed.

Another important reason why we first allocate the array object into *S* is because the array is often much larger than the top and tuple objects. It is much more efficient to allocate it directly into the space it belongs to rather than allocating it somewhere else and moving it later.

Table 1 shows our allocation policies for different types of objects in an RDD. For RDDs with tag “DRAM”, array objects are allocated directly into the DRAM component of the old generation if it has enough space. Otherwise, they have to be allocated in the NVM component. For RDDs with tag “NVM”, array objects are allocated directly into the NVM component. Intermediate RDDs without tags are all allocated in the young generation (DRAM). Most of them end up dying there and never get promoted, while a small number of objects that eventually become old enough will be promoted to the NVM space of the old generation. Top RDD objects and data tuple objects, as discussed earlier, are all allocated into the young generation and moved later by the GC to the spaces containing their corresponding arrays.

4.2 Implementation and Optimization

This subsection describes our implementation techniques and various optimizations.

Table 1. Panthera’s Allocation Policies

Tag	Obj Type	Initial Space	Final Space
DRAM	RDD Top	Young Gen.	DRAM of Old Gen.
	RDD Array	DRAM of Old Gen.	DRAM of Old Gen.
	Data Objs	Young Gen.	DRAM of Old Gen.
NVM	RDD Top	Young Gen.	NVM of Old Gen.
	RDD Array	NVM of Old Gen.	NVM of Old Gen.
	Data Objs	Young Gen.	NVM of Old Gen.
NONE	RDD Top	Young Gen.	Young Gen. or NVM of Old Gen.
	RDD Array	Young Gen.	Young Gen. or NVM of Old Gen.
	Data Objs	Young Gen.	Young Gen. or NVM of Old Gen.

4.2.1 Passing Tags. Right before each materialization point (i.e., the invocation of `persist` or a Spark action), our analysis inserts a call to a native method `rdd_indicator(rdd, tag)`, with the RDD’s top object (`rdd`) and the inferred memory tag (`tag`) as the arguments. This method first sets a thread-local state variable to DRAM or NVM, according to the tag, informing the current thread that a large array for an RDD will be allocated soon. Next, `rdd_indicator` sets the `MEMORY_BITS` of the top object `rdd` based on `tag`. Regardless of where it currently is, this top object will eventually be moved by the GC to the space corresponding to `tag`.

The thread then transitions into a “wait” state, waiting for this large array. In this state, the first allocation request for an array whose length exceeds a user-defined threshold (i.e., a million used in our experiments) is recognized as the RDD array. Panthera then allocates the array directly into the space indicated by `tag`. To implement this, we modified both the fast allocation path, assembly code generated by the JIT compiler, and the slow path, functions implemented in C++. After this allocation, the state variable is reset and the thread exits the wait state. If `tag` is null, the array is allocated in the young generation, preferably through the **thread-local allocation buffer (TLAB)**, and the `MEMORY_BITS` of the top object remains as the default value (00).

4.2.2 Object Migration. There are two major challenges in how to move objects: *cross-generation migration* and *object compaction*. As Panthera piggybacks on a generational GC where a minor GC is triggered when JVM is unable to allocate space for a new object, objects in the young generation that survive several minor GCs are deemed “long-lived” and moved into the old generation. The major GC is triggered when the old generation is full. We leverage this opportunity to move together objects that belong to the same logical RDD—as discussed earlier, these objects might not have been allocated in the same space initially.

Minor GC. To do this, we modified the minor collection algorithm in the Parallel Scavenge GC on which Panthera is built. The existing minor GC contains three tasks: root-task, which performs object tracing from the roots (e.g., stack and global variables); old-to-young-task, which scans references from objects in the old generation to those in the young generation to identify (directly or transitively) reachable objects; and steal-task, which performs work stealing for load balancing. To support our object migration, we split old-to-young-task into a DRAM-to-young-task and NVM-to-young-task, which find objects that should be moved into the DRAM and NVM parts of the old generation, respectively.

For these two tasks, we modified the tracing algorithm to propagate the tag—for example, scanning a reference from a DRAM-based RDD array (with tag “DRAM”) to a tuple object (in the young generation) propagates the tag to the tuple object (by setting its `MEMORY_BITS`). Hence, when tracing is done, all objects reachable from the array have their `MEMORY_BITS` set to the same value as that of

the array. In the original GC algorithm, an object does not get promoted from the young to the old generation until it survives several minor GCs (in this article, we use the threshold of 15). In Panthera, however, we move the objects whose `MEMORY_BITS` is set as 01 (10) in tracing immediately to DRAM (NVM) space in the old generation. We refer to this mechanism as *eager promotion*. Objects whose `MEMORY_BITS` is not set, 00, in tracing belong to intermediate RDDs or are control objects not associated with any RDDs. The migration of these objects follows the original algorithm, that is, they will be moved only if they survive several minor GCs. During the eager promotion, if there is a lack of free DRAM space for the old generation, Panthera will put the corresponding objects into NVM and let major GC adjust the data layout during execution.

Furthermore, we also need to move RDD top objects to the appropriate part of the old generation. These top objects, whose `MEMORY_BITS` was set by the instrumented call to `rdd_indicator` at their materialization points, are visited when root-task is executed because these objects are referenced directly by stack variables. We modified the root-task algorithm to identify objects with the set `MEMORY_BITS`. These RDD top objects will also be moved to (the DRAM (01) or NVM (10) space of) the old generation by the minor GC.

Major GC. When a major GC runs, it performs memory compaction by moving objects together (in the old generation) to reduce fragmentation and improve locality. We modified the major GC to guarantee that compaction does not occur across the boundary between DRAM and NVM. Furthermore, when the major GC performs a full-heap scan, Panthera re-assesses, for each RDD array object, where the object should actually be placed based on the RDD's runtime access frequency. This frequency is measured by counting, using instrumentation, how many times a method (e.g., `map` or `reduce`) has been invoked on this RDD object. The RDDs are ranked based on the access frequency. The most frequently accessed RDDs will be migrated to DRAM if they are misplaced in NVM. If there isn't enough DRAM space for these high-ranking RDDs, Panthera will evict the RDDs with lower access frequencies from the DRAM.

We maintain a hash table that maps each RDD object to the number of calls made on the object. Our static analysis inserts, at each such call site, a JNI (Java Native Interface) call that invokes a native JVM method to increment the call frequency for the RDD object. Frequently (infrequently) accessed array objects are moved from the NVM (DRAM) space to the DRAM (NVM) space within the old generation and all objects reachable from these arrays are moved as well. Their `MEMORY_BITS` will be updated accordingly. At the end of each major GC, the frequency for each RDD is reset.

The DRAM space of the old generation can be quickly filled up as it is much smaller than the NVM space. When the DRAM space is full, the minor GC moves all objects from the young generation to the NVM space of the old generation regardless of their memory tags.

Conflicts. If an object is reachable from multiple references and different tags are propagated through them, a conflict occurs. As discussed earlier, we resolve conflicts by giving "DRAM" higher priority than "NVM". As long as the object receives "DRAM" from any reference, it is a DRAM object and will be moved to the DRAM space of the old generation.

4.2.3 Card Optimization. In OpenJDK, the heap is divided into many *cards*, each representing a region of 512 bytes. Every object can take one or more cards, and the write barrier maintains a card table that marks certain cards dirty upon reference writes. The card table can be used to efficiently identify references during tracing. For example, upon `a.f = b`, the card that contains the object referenced by `a` is set to dirty. When a minor GC runs, the old-to-young scavenge task cleans a card if the target objects of the (old-to-young) references contained in the memory region represented by the card have been copied to the old generation.

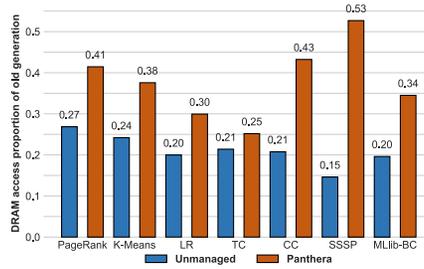


Fig. 6. DRAM access proportion in old generation.

However, if a card contains two large arrays (say A and B)—e.g., A ends in the middle of the card while B starts there immediately—significant inefficiencies can result when they are scanned by two different GC threads. The card would remain dirty even if all objects referenced by A and B have been moved from the young to the old generation—neither thread could clean the card due to its unawareness of the status of the array scanned by another thread. This would cause every minor GC to scan every element of each array in the dirty card until a major GC occurs.

This is a serious problem for Big Data applications that make heavy use of large arrays. Shared cards exist pervasively when these arrays are frequently allocated and deallocated. Frequent scanning of such cards with multiple threads can incur a large overhead on NVM due to its higher read latency and reduced bandwidth. We implemented a simple optimization that adds an *alignment padding* for the allocation of each RDD array to make the end of the array align with the end of a card. Although this leads to space inefficiencies, the amount of wasted space is small (e.g., less than 512 bytes for each array of hundreds of megabytes) while card sharing among arrays is completely eliminated, resulting in substantial reduction in GC time.

5 PROFILING-GUIDED OPTIMIZATION

As described in Section 3, the core idea behind Panthera is to statically infer RDDs’ memory tags and pass them to the runtime system to instruct objects migration. In order to analyze the runtime behavior of the objects more precisely, we propose the **profiling-guided optimization (PGO)** in this section.

5.1 Memory Access Distribution

Figure 6 shows the proportion of DRAM memory accesses to total memory accesses in old generation. Compared with Unmanaged, Panthera increased the proportion of DRAM access by 16.7% on average, which demonstrates that Panthera’s strategy can effectively increase DRAM access proportion. In this section, we design a new experiment to analyze the access behaviors of objects in the old generation with a fine granularity.

We divide the old generation space into chunks of 1-GB size, and plot the access behaviors for each chunk in Figure 7, using four applications, i.e., LR, TC, GraphX-CC, and MLlib. In particular, the horizontal axis uses Chunk ID to represent all the chunks, and the vertical axis represents the access proportion for each chunk, i.e., the number of accesses on the chunk divided by the number of accesses on all chunks. The chunks in Figure 7 are divided into four categories, (1) red, frequently accessed chunks and placed in DRAM, (2) yellow, infrequently accessed chunks but placed in DRAM, (3) blue, frequently accessed chunks but placed in NVM, and (4) green, infrequently accessed chunks and placed in NVM. The red and green dots show the chunks that are correctly recognized and placed. The yellow and blue dots show the chunks that are wrongly placed.

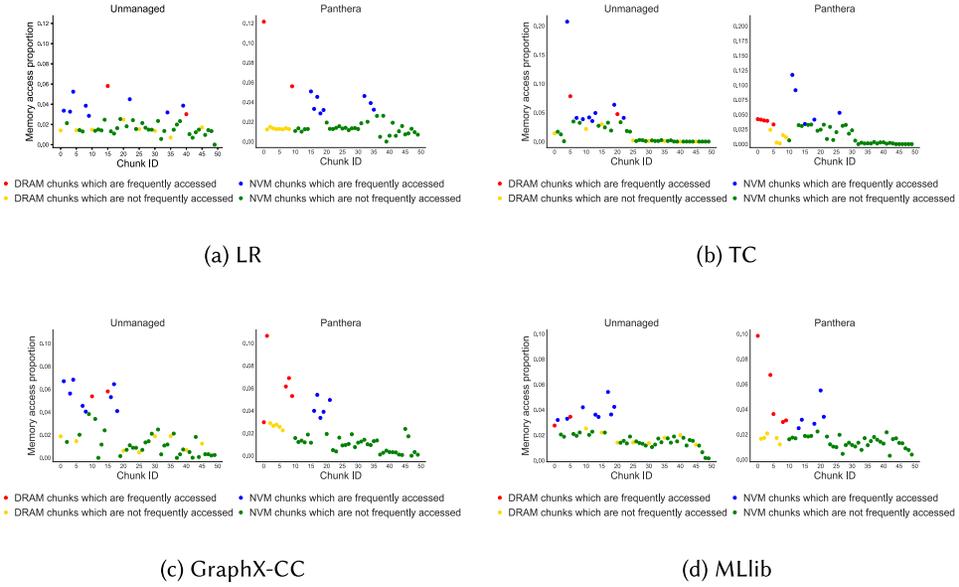


Fig. 7. Spatial distribution of memory access.

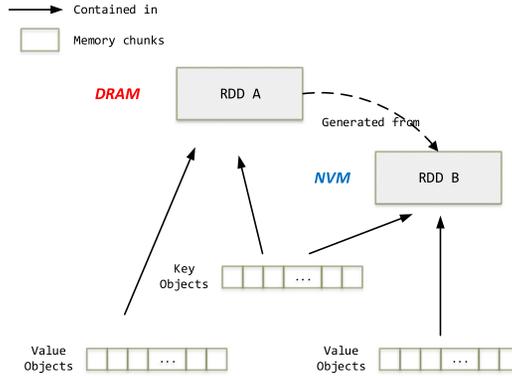


Fig. 8. Example of data sharing.

From the figure, we can see that, comparing with Unmanaged, Panthera can place more frequently accessed chunks on DRAM and infrequently accessed chunks on NVM. Take the MLib (Figure 7(d)) for example, among the top-10 frequently-accessed chunks, the Unmanaged approach allocates two chunks on DRAM, while Panthera allocates five chunks on DRAM. However, there still exist some frequently accessed chunks that are placed on NVM, as shown by the blue dots, and some infrequently accessed chunks that are placed on DRAM, as shown by the yellow dots, which are undesirable. The reason of the data misplacement is that Panthera without PGO treats the RDD as a whole and can't distinguish the access frequency difference within the RDDs.

In summary, our key finding on Spark is that the data belonging to the same RDD do not always exhibit similar access patterns. This finding motivates us to introduce finer-grained chunk placement decision into Panthera. Therefore, we only need to perform chunk-based profiling to the data that belonging to the annotated RDD, rather than profiling for all data objects during the program execution. Therefore, we integrate the coarse-grained RDD-level analysis globally and

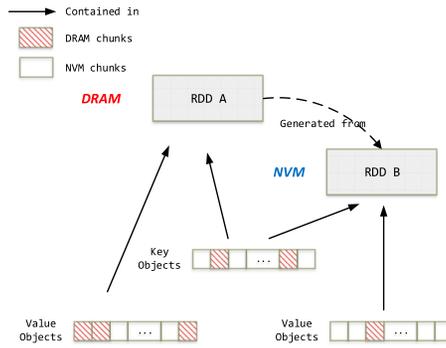


Fig. 9. Memory tags for each chunk after the profile-guided optimization.

fine-grained chunk-based profiling for some special RDDs. This would obtain precise memory access patterns in a lightweight manner.

5.2 Opportunity for Runtime Optimization

As RDDs are multi-layer Java data structures, and Panthera without PGO analyzes the behavior at the granularity of RDD, thus all objects belonging to one RDD would be identified to have the same access pattern and lifetime. However, some objects might be shared by multiple RDDs that are identified to have different behaviors, and it brings an opportunity for more precisely allocating the objects across DRAM/NVM.

As discussed in Section 3.1, we discussed an example of object sharing, and statically, we proposed an algorithm and attempted to assign the same memory tag to the RDDs sharing the objects as shown in Figure 3. However, the attempt might fail, and there might bring conflicts when inferring the tag for the objects from different RDDs. For example, in Figure 8, RDD B is generated from RDD A, B and A receive the tag of “NVM” and “DRAM”, respectively. Therefore, when B is generated, the data objects that are shared by A and B will be moved into “NVM”, even if some of them belong to A which is tagged “DRAM”.

Based on the observation, we have the opportunity to dynamically switch the memory tag for the shared objects, i.e., using the tag of “NVM” when accessing B, and “DRAM” when accessing A. For this purpose, we need to refine the analysis, from the granularity of RDDs to chunks, and leverage the runtime behaviors of the objects.

5.3 Optimization

To characterize the behaviors at the granularity of chunks rather than RDDs and allocate more frequently accessed chunks to DRAM, we propose a profile-guided approach, which works as follows: First, we collect the memory access distribution for all the chunks by using VTune, with the chunk size of CS , as shown in Figure 7. Second, from the access distribution, we select the top- K frequently accessed chunks, where K is determined by the DRAM capacity of old generation C and the chunk size CS using the equation of $K = C/CS$. Finally, the IDs for the K chunks are passed to the JVM, and when JVM starts, we bind the top- K chunks to “DRAM”, and other chunks to “NVM”, by invoking the `mbind` system call. to determine the memory address for each chunk. For example, we assume chunk size is 1GB. When a 64-GB heap is used and DRAM to memory ratio is 1/3 and the nursery space is 1/6 of the heap size, there is 10.66-GB DRAM in old generation. Then, the IDs for the top 11 frequently accessed chunks are passed to the JVM where the top-10 chunks and the first 0.66 GB of the 11th chunk will be bound into “DRAM”.

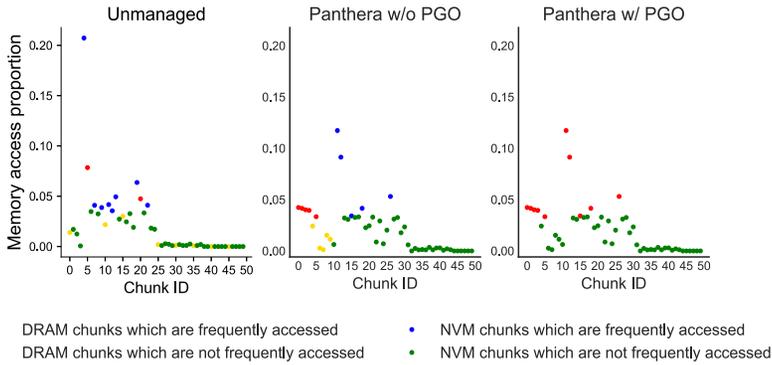


Fig. 10. Allocation results of DRAM chunks after the profile-guided optimization.

Figure 9 shows the enhanced Panthera framework with the profile-guided optimization, for executing the example in Figure 8. In particular, in Spark analyzer, the RDDs B and A are tagged as “DRAM” and “NVM”, respectively. When JVM starts, the profiled K chunks would be bound to DRAM, and other chunks to NVM. Therefore, from the perspective of RDD, B would be still allocated to the DRAM part of the old generation, and A would be still allocated to the NVM part of the old generation. However, with our underlying chunk binding, only the profiled hot K chunks of A and B would be actually allocated to DRAM.

Therefore, we leverage the profiled access frequency of the memory chunks, and refine the statically determined DRAM/NVM partition, to allocate only the real hot chunks on DRAM. Figure 10 shows the results of applying the PGO to Panthera. From this figure, we see that, for the Panthera with PGO, most of the frequently accessed chunks are correctly marked with “DRAM” tags and placed in DRAM. Compared to the Unmanaged and Panthera without PGO, Panthera with PGO significantly improved the accumulated access proportion on DRAM of Old generation from 19.55% and 25.2% to 53.5%.

6 DISCUSSION ON APPLICABILITY AND GENERALITY

Panthera’s design includes two key individual mechanisms, i.e., analyzing the data access patterns which is framework-dependent such as the Spark Analyzer and QuickCached Analyzer in Section 3, and a set of framework-independent APIs that makes pretenuring, migration, and dynamic monitoring easy for any in-memory big data system using large arrays as backbone data structures. In our design, the clear and predictable data access patterns are connections between different on-top frameworks and the enhanced JVM.

6.1 General Memory Management Policies

The enhanced JVM is general since the data placement and migration mechanism provided by the Panthera runtime system can be employed to manage memory for any Big Data systems that have clear and predictable data access patterns. Examples include Apache Hadoop, Apache Flink, or database systems such as Apache Cassandra.

Panthera determines the data placement and migration via three general policies. These policies are integrated into the enhanced JVM thus they are generally applicable to other JVM-based big data systems.

- First, some data structures can be pre-tenured with some tags according to their behaviors that can be statically determined. Panthera would allocate these data structures directly into the space indicated by the corresponding tags.

- Second, some data structures are required to collect their behaviors and determine their placement at runtime. In our current Panthera implementation, the placement of these data structures depends on their access frequencies and lifetimes, thus Panthera would collect these runtime characteristics and make the placement decision during the application execution. Furthermore, Panthera can be extended to integrate new memory access characteristics and new policies.
- Third, some data structures are required to make finer-grained placement decision. For these data structures, Panthera leverages a finer-grained chunk-based memory access profiling approach that enables placements of different chunk at different memory spaces.

However, Panthera currently does not support applications whose data access patterns cannot be distinguished clearly, e.g., applications with random memory access.

6.2 Framework-Specific Access Pattern Annotations/Analyzers

The memory access patterns are obtained via user annotations together with static analyzers.

Panthera provides two major APIs, one for pre-tenuring data structures with tags and a second for dynamic monitoring and migration. The first API takes as input an array and a tag, performing data placement as discussed earlier in Section 4. The tag can come from the developer's annotations in the program or from a static analysis that is designed specifically for the framework to be optimized.

To illustrate, consider Apache Hadoop where both a map worker and a reduce worker may need to hold large data structures in memory. Some of these data structures are loaded from HDFS as immutable input, while others are frequently accessed. In the case of HashJoin, which is a building block for SQL engines, one input table is loaded entirely in memory while the second table is partitioned across map workers. If map workers are executed in separate threads, they all share the first table and join their own partitions of the second table with it. The first table is long-lived and frequently accessed. Hence, it should be tagged DRAM and placed in the DRAM space of the old generation, while different partitions of the second table can be placed in the young generation and they will die there quickly.

Panthera's second API takes as input a data structure object to track the number of calls made on the object. If this API is used to track the access frequency of the data structure, the data structure (and all objects reachable from it) would not be pre-tenured (as specified by the first API), but rather, they are subject to dynamic migration performed in the major GC. We can use this API to dynamically monitor certain objects and migrate them if their access patterns are not easy to predict statically.

Use of these above two APIs enables a flexible allocation/migration mechanism that allows certain parts of the data structure (e.g., for which memory tags can be easily inferred) to be pre-tenured and other parts to be dynamically migrated. Furthermore, the framework-specific memory access pattern analyzers would perform def-use analysis to propagate the user annotations to the runtime system. In Section 3, we demonstrated two individual analyzers for Spark and QuickCached, respectively, which are not easily reusable to other new frameworks.

6.3 Apply Panthera to a New Framework

When extending Panthera to a new framework, we need to consider the following key issues:

First, we might need to introduce new framework-specific annotations so that the runtime system would have the knowledge about the key data structures. For example, as described in Section 3, to apply Panthera to QuickCached, the `@CoreHashObject` annotation is introduced to illustrate that this is the data structure for the global hash table in QuickCached. In particular, the

annotations are designed together with the static analyzers and deliver application-level knowledge to the analyzers.

Second, we need to design a new framework-specific static analyzer to expose the memory access patterns. For example, a new static QuickCached analyzer is developed to identify the core data structure based on user's annotations and allocate an auxiliary array. The insight behind is that only a small fraction of the hashtable is frequently accessed, and the insight guides us to introduce the auxiliary array to hold these frequently accessed data. Meanwhile, the auxiliary array would be tagged by the analyzer so that it can be pre-promoted into the DRAM space of the old generation. Note that the analyzers can leverage user annotations and some framework-specific heuristics to obtain more precise memory access patterns.

With the framework-specific annotations and analyzers, the underlying runtime mechanism would leverage the second APIs and dynamically determine the data structure placement and migrations, without any framework-specific modifications. For example, for QuickCached, the runtime system would collect the frequently accessed data to the auxiliary array, thus these data would be allocated on DRAM.

7 EVALUATION

We have added/modified 9,186 lines of C++ code in OpenJDK (build jdk8u76-b02) to implement the Panthera GC, and written 979 lines of Scala code to implement the static analysis for Spark and 762 lines of Java code for QuickCached.

7.1 NVM Emulation and Hardware Platform

Most of the prior works on hybrid memories used simulators for experiments. However, none of them support Java applications well. We cannot execute managed-runtime-based distributed systems on these simulators. There also exist emulators such as Quartz [75] and P MEP [29] that support emulation of NVM for large programs using commodity multi-socket (NUMA) hardware, but neither Quartz nor P MEP could run OpenJDK. These emulators require developers to use their own libraries for NVM allocation, making it impossible for the Panthera GC to migrate objects without re-implementing the entire allocator and GC from scratch using these libraries.

As observed in [10] and [75], NUMA's remote memory latency is close to NVM's latency, and hence, researchers have used a NUMA architecture as the baseline to measure emulation accuracy. Following this observation, we built our own emulator on NUMA machines to emulate hybrid memories for JVM-based Big Data systems.

We followed Quartz [75] when implementing our emulator. Quartz has two major components: (1) it uses the *thermal control register* to limit the DRAM bandwidth; and (2) it creates a daemon thread for each application process and inserts delay instructions to emulate the NVM latency. For example, if an application's CPU stall time is S , Quartz scales the CPU stall time to $S \times \frac{NVM_latency}{DRAM_latency}$ to emulate the latency effect of NVM. For (1), we used the same thermal control register to limit the read/write bandwidth. Like Quartz, we currently do not support different bandwidths for reads and writes. For (2), we followed Quartz's observation to use the latency of NUMA's remote memory to model NVM's latency.

An alternative approach to emulating NVM's latency is to instrument loads/stores during JIT compilation, injecting a software-created delay at each load/store. The limitation of this approach, however, is that it does not account for caching effects and memory-level parallelism.

We used one CPU to run all the computation, the memory local to the CPU as DRAM, and the remote memory as NVM. In particular, DRAM and NVM are emulated, respectively, using two local and two remote memory channels. The performance specifications of the emulated NVM are

Table 2. Emulated DRAM and NVM Parameters

	DRAM	NVM
Read latency (ns)	120	300
Bandwidth (GB/s)	30	10 (limited by the thermal control register)
Capacity per CPU	100s of GBs	Terabytes
Estimated price	5×	1×

the same as those used in [75], reported in Table 2. To emulate NVM’s slow write speed, we used the thermal control register to limit the bandwidth of remote memory—the emulated NVM is full duplex with 10 GB/s for read and write bandwidth each. The remote memory’s latency in our setting is 2.5× of that of the local memory.

Energy Estimation. We followed Lee et al. [49] to estimate energy for NVM. We used Micron’s DDR4 device specifications [61] to model DRAM’s power. NVM’s energy has a *static* and *dynamic* component. The static component is negligible compared with DRAM [50]. The dynamic component consists of the energy consumed by reads and writes. PCM array reads consume about 2.1× larger energy than DRAM due to its need for high temperature operation [49].

NVM writes consume much more energy than DRAM writes. Upon a row-buffer miss, the energy consumed by each write has three components: (1) an *array write* that evicts data from the row buffer into the bank array, (2) an *array read* that fetches data from the bank array to the row buffer, and (3) a *row buffer write* that writes new data from the CPU last level cache to the row buffer. Assuming the row-buffer miss ratio is 0.5, we computed these three components separately by considering the row buffer’s write energy (1.02 pJ/bit), size (i.e., 8K bits for DRAM [61], 32-bit-wide partial writeback to NVM [49]) and miss rate (0.5), as well as the array’s write-back energy (16.8pJ/bit × 7.6% for NVM) and read energy (2.47pJ/bit for NVM). The factor of 7.6% is due to Lee et al.’s optimization [49] that writes only 7.6% of the dirty words back to the NVM array.

CPU’s uncore events, collected with *VTune* [7], were employed to compute the numbers of reads and writes. In particular, the events we used were `UNC_M_CAS_COUNT.RD` and `UNC_M_CAS_COUNT.WR`. *VTune* can also distinguish reads and writes from/to local and remote memories.

7.2 Experiment Setup

We set up a small cluster to run Spark with one master node and one slave node—these two servers have a special Intel chipset with a “scalable memory buffer” that can be tuned to produce the 2.5× latency for remote memory accesses, which matches NVM’s read/write latency. Since our focus is *not* on distributed computing, this cluster is sufficient for us to execute real workloads on Spark and understand their performance over hybrid memories. Table 3 reports the hardware configurations of the Spark master and Spark slave nodes. Each node has two 8-core CPU and the Parallel Scavenge collector on which Panthera was built creates 16 GC threads in each GC to perform parallel tracing and compaction.

The negative impact of the GC latency increases with the number of compute nodes. As reported in [57], a GC run on a single node can hold up the entire cluster—when a node requests a data partition from another server that is running GC, the requesting node cannot do anything until the GC is done on the second node. Since Panthera can significantly improve the GC performance on NVM, we expect Panthera to provide even greater benefit when Spark is executed on a large NVM cluster.

Table 3. Hardware Configuration for Our Servers

Arch	NUMA, 4 sockets QPI 6.4 GT/S, directory-based MESIF
CPU	E7-4809 v3 2.00 GHz, 8 cores, 16 HW threads
L1-I	8 way, 32 KB/core, private
L1-D	8 way, 32 KB/core, private
L2	8 way, 256 KB/core, private
L3	20 way, 20 MB, shared
Memory	DDR 4, 1,867 MHz, SMI 2 channels

System Configurations. Each CPU has a 128-GB DRAM. We reserved 8 GB of DRAM for the OS and the maximum heap amount of DRAM that can be used for Spark is 120 GB. We experimented with two different heap sizes for the Spark-running JVM (64 GB and 120 GB) and three different DRAM sizes (1/4, 1/3, and 100% of the heap size; the rest of the heap is NVM). For QuickCached, we experimented a 64-GB heap size with two different DRAM sizes (1/3 and 100% of the heap size). The configuration with 100% DRAM was used as a baseline to compute the overhead of Panthera under hybrid memories.

Prior works on NVM often used smaller DRAM ratios in their configurations. For example, Write Rationing [11] used 1-GB DRAM and 32-GB NVM in their experiments. However, as we deal with Big Data systems, it would not be possible for us to use a very small DRAM ratio—in our experiments, a regular RDD consumes 10-30-GB memory, and hence, we had to make DRAM large enough to hold at least one RDD.

The nursery space is placed entirely in DRAM. We have experimented with several different sizes (1/4, 1/5, 1/6, and 1/7 of the heap size) for the nursery space. The performance differences between the 1/4, 1/5, and 1/6 configurations were marginal (even under the original JVM), while the configuration of 1/7 led to worse performance. We ended up using 1/6 in our experiments for both Spark and QuickCached to achieve good nursery performance and simultaneously leave more DRAM to the old generation.

Programs and Datasets. For Spark, We selected a diverse set of seven programs. Table 4 lists these programs, the datasets used to run them and their memory footprints. These are representative programs for a wide variety of tasks including data mining, machine learning, graph and text analytics. PR, KM, LR, and TC run directly on Spark; CC and SSSP are graph programs running on GraphX [33], which is a distributed graph engine built over Spark; BC is a program in MLlib, a machine learning library built on top of Spark. We used real-world datasets to run all the seven programs. Note that although the sizes of these input datasets are not very large, there can be large amounts of intermediate data generated during the computation.

To evaluate the performance of QuickCached, we use the **Yahoo! Cloud Serving Benchmark (YCSB)** [22]. YCSB is a benchmark suite commonly used to evaluate the performance of cloud storage services. We run its A, B, C, D, and F workloads after loading the databases with 30 million records. Each record is 1 KB by default. For each workload, we perform 10 million operations.

Baselines. Our initial goal was to compare Panthera with both Espresso [80] and Write Rationing [11]. However, neither of them is publicly available. Espresso proposes a programming model for developers to develop new applications. Applying it to Big Data systems would mean that we need to rewrite each allocation site, which is clearly not practical. In addition, Espresso does not migrate objects based on their access patterns.

Table 4. Spark Programs, Datasets and Memory Footprints

Program	Dataset	Initial Size	Footprint
PageRank (PR)	Wikipedia Full Dump, German [4]	1.2 GB	63.0 GB
K-Means (KM)	Wikipedia Full Dump, English [4]	5.7 GB	49.5 GB
Logistic Regression (LR)	Wikipedia Full Dump, English [4]	5.7 GB	63.4 GB
Transitive Closure (TC)	Notre Dame Webgraph [2]	21 MB	43.8 GB
GraphX-Connected Components (CC)	Wikipedia Full Dump, English [4]	5.7 GB	59.1 GB
GraphX-Single Source Shortest Path (SSSP)	Wikipedia Full Dump, English [4]	5.7 GB	62.2 GB
Mllib-Naive Bayes Classifiers (BC)	KDD 2012 [1]	10.1 GB	63.1 GB

The Write Rationing GC has two implementations: *Kingsguard-Nursery* (KN) and *Kingsguard-Writes* (NW). KN places the young generation in DRAM and the old generation in NVM. KW also places the young generation in DRAM. Different from KN, KW monitors object writes and dynamically migrates write-intensive objects into DRAM. Although we could not directly compare Panthera with these two GCs, we have implemented similar algorithms in OpenJDK. Under KW, almost all persisted RDDs were quickly moved to NVM. The frequent NVM reads from these RDDs, together with write barriers used to monitor object writes, incurred an average of 41% performance overhead for our benchmarks. This is because Big Data applications exhibit different characteristics from regular, non-data-intensive Java applications.

KN appears to be a good baseline at the first sight. However, implementing it naïvely in the Parallel Scavenge collector can lead to non-trivial overhead—the reduced bandwidth in NVM can create a huge impact on the performance of a multi-threaded program; this is especially the case for Parallel Scavenge that attempts to fully utilize the CPU resources to perform parallel object scanning and compaction.

To obtain a better baseline, we placed the young generation in DRAM and supported the old generation with a mix of DRAM and NVM. In particular, we divided the virtual address space of the old generation into a number of chunks, each with 1 GB, and used a probability to determine whether a chunk should be mapped to DRAM or NVM. The probability is derived from the DRAM ratio in the system. For example, in a system where the DRAM-to-memory ratio is 1/4 (1/4 DRAM), each chunk is mapped to DRAM with 1/4 probability and to NVM with 3/4 probability. Note that this is common practice [32, 77] to utilize the combined bandwidth of DRAM and NVM. We refer to this configuration as *unmanaged*, which outperforms both KN and KW for our benchmarks.

There are also some OS-level-based data migration works, such as Thermostat [9], Translation Ranger [86], and HeteroOS [43]. We tried to port these frameworks to our emulated NVM platform, but none of these works fit for the benchmarks we used. For example, when running on Thermostat, the Spark applications always get stuck during the execution and the HeteroOS targets at the hybrid memory in virtualized environments instead of running on the bare-metal machines as Panthera does. Translation Ranger targets at speeding up the virtual-to-physical memory address translation by actively coalescing fragmented pages. It's an orthogonal optimization to Panthera, so we didn't include it in the evaluations. In order to evaluate the OS-level hybrid memory management policy, we utilize the kernel **LRU (Least Recently Used)** based paging system to do the data migration management. We created a ramdisk on the emulated NVM and mount it as the swap partition. We tuned the performance of the paging system to the best according to the

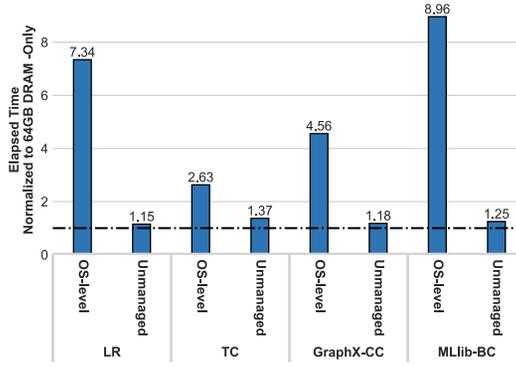


Fig. 11. Overall performance comparison between OS-level management and *unmanaged*. The heap size is 64GB and DRAM to memory ratio is 1/3.

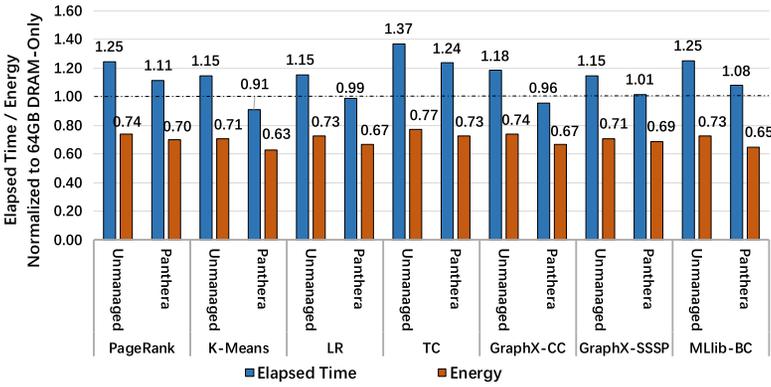


Fig. 12. Overall performance and energy results of Spark under a 64-GB heap; DRAM to memory ratio is 1/3.

state-of-art works [12, 59]. Under this settings, the NVM works as a secondary memory, similar with the Optane DC Memory Mode [37]. The kernel evicts the least recently used data to NVM and keeps the most recently used data in DRAM. As Figure 11 shows, the OS-level management policy is much worse than the *unmanaged*. Compared to *unmanaged*, the slowdown of OS-level management can reach to up 6.20 \times . This is because the GC always messes up the data layout placed by the OS, which cause much more useless data migration overhead, as we described in the Section 1. Hence, we utilize the *unmanaged* as baseline in the subsequent evaluations.

7.3 Performance and Energy of Panthera without PGO

Figure 12 reports the overall performance and energy results of Spark when a 64-GB heap is used and then DRAM-to-memory ratio is 1/3 (1/3 DRAM). The performance and energy results of each configuration are normalized w.r.t. those of the 64-GB DRAM-only version. The energy results in our experiments include the energy consumption of Panthera runtime, but do not include the energy consumption of the static analyzer and profiling tools. Compared to the DRAM-only version, the unmanaged version reduces energy by 26.7% with a 21.4% execution time overhead. In contrast, Panthera reduces energy by 32.3% at a 4.3% execution time overhead.

When the heap size is 120 GB (not shown in Figure 12, but summarized later in Figure 14 and Figure 15), the unmanaged version reduces energy by 39.7% at a 19.3% execution time overhead. In contrast, Panthera reduces energy by 47.0% with less than 1% execution time overhead. Clearly,

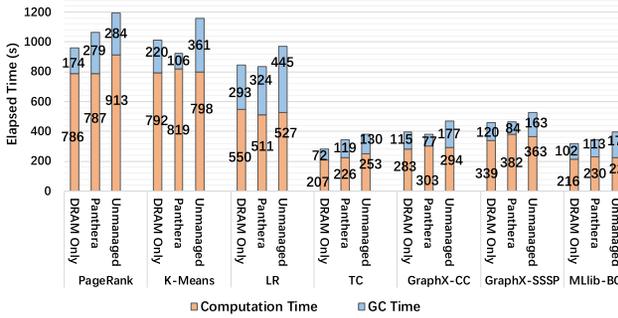


Fig. 13. GC performance (64-GB heap).

considering the RDD semantics in data placement provides significant benefits in both energy and performance.

GC Performance. To understand the GC performance, we broke down the running time of each program into the mutator and GC time; these results (under the 64-GB heap) are shown in Figure 13. Compared to the baseline, the unmanaged version introduces performance overhead of 60.4% and 6.9% in the GC and computation, respectively; while for Panthera these two overheads are, respectively, 4.7% and 4.5%. Under the 120-GB heap, the GC performance overhead of the unmanaged version and Panthera are, respectively, 58.0% and 3.1%. Note that, due to large amounts of intermediate data generated, the GC is frequently triggered for these programs.

Since the GC is a memory-intensive workload, inappropriate data placement can lead to significantly higher memory access time and thus a large penalty. The penalty comes from two major sources. First, NVM’s limited bandwidth (which is about 1/3 of that of DRAM) has a large negative impact on the performance of Parallel Scavenge, which launches 16 threads to perform parallel tracing and object copying in each (nursery and full-heap) GC. Given this high degree of parallelism, the performance of the nursery GC is degraded significantly when scanning objects in NVM. Second, object tracing is a read-intensive task, which suffers badly from NVM’s higher read latency.

Panthera improves the GC performance by pretenuring frequently accessed RDD objects in DRAM and performing optimizations including *eager promotion* (Section 4.2.2) and *card padding* (Section 4.2.3). *Eager promotion* reduces the cost of (old-to-young) tracing in each minor GC, while *card padding* eliminates unnecessary array scans in NVM, which are sensitive to both latency and bandwidth. A further breakdown shows that *eager promotion*, alone, contributes an average of 9% of the total GC performance improvement. The contribution of *card padding* is much more significant—without this optimization, the GC time increases by 60% due to the impact of NVM’s substantially limited bandwidth and increased latency on the performance of parallel card scanning. In fact, this impact is so large that the other optimizations would not work well when card padding is disabled.

Varying Heaps and Ratios. To understand the impact of the heap sizes and DRAM ratios (DRAM to total memory), we have conducted experiments with two heap sizes (64 GB, 120GB) and two DRAM ratios (1/3, 1/4) on four programs PR, LR, CC, and BC. Figure 14 reports the time results of these configurations. Panthera’s time overheads are, on average, 9.5%, 3.4%, 2.1%, and 0%, respectively, under the four configurations (64 GB, 1/4), (64 GB, 1/3), (120 GB, 1/4), and (120 GB, 1/3). The overheads for the unmanaged version are 25.9%, 20.9%, 23.9%, and 19.3%, respectively, under these same four configurations.

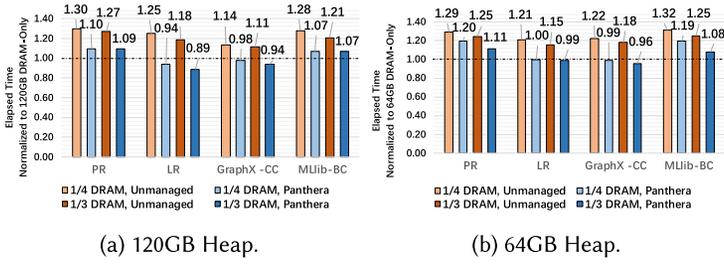


Fig. 14. Performance for two DRAM ratios + two heaps.

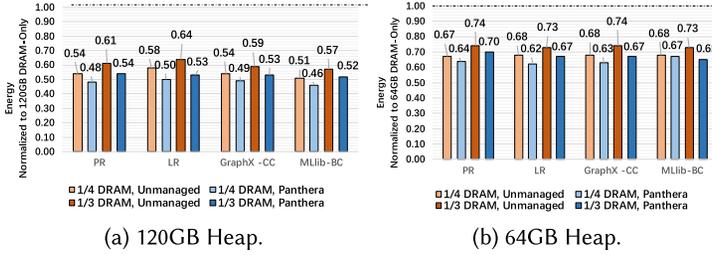


Fig. 15. Energy for two DRAM ratios + two heaps.

We make two interesting observations. First, Panthera is more sensitive to the DRAM ratio than the heap size. The time overhead can be reduced by almost 10% when the DRAM ratio increases from 1/4 to 1/3. The reason is that more frequently accessed RDDs are moved to DRAM, reducing the memory latency and bandwidth bound of NVM. Another observation is that the unmanaged version is much less sensitive to DRAM ratio—the time overhead is reduced by only 5% when the DRAM ratio increases to 1/3. This is because arbitrary data placement leaves much of the frequently accessed data in NVM, making CPUs stall heavily when accessing NVM.

Figure 15 depicts the energy results for the two heaps and two DRAM/NVM ratios. For the 64-GB heap, the unmanaged version reduces energy by an average of 32.2% and 26.5%, respectively, under the 1/4 and 1/3 DRAM ratio, while Panthera reduces energy by 36.0% and 32.7% under these same ratios. The energy reductions for the 120-GB heap are much more significant—the unmanaged version reduces energy by 45.7% and 39.7%, respectively, under the 1/4 and 1/3 DRAM ratios, while the energy reduction under Panthera increases to 51.7% and 47.0% for these two ratios.

We also evaluate the prices of NVM and DRAM to show the hardware cost savings that benefit from using the hybrid memory. The price of the cheapest NVM is \$7.85 per GB and the cheapest DRAM is \$16.61 per GB [35]. Compared with DRAM-only, using the hybrid memories with DRAM ratio 1/3 can reduce 35.2% hardware costs, and even reduce 39.6% when with DRAM ratio 1/4.

Results for QuickCached. Figure 16 reports the overall performance and energy results of QuickCached when a 64-GB heap is used and DRAM to memory ratio is 1/3 (1/3 DRAM). The performance and energy results of each configuration are normalized w.r.t. those of the 64-GB DRAM-only version. Compared to the DRAM-only version, the unmanaged version reduces energy by 25.0% with a 9.1% execution time overhead. In contrast, Panthera reduces energy by 28.7% at a 5.2% execution time overhead.

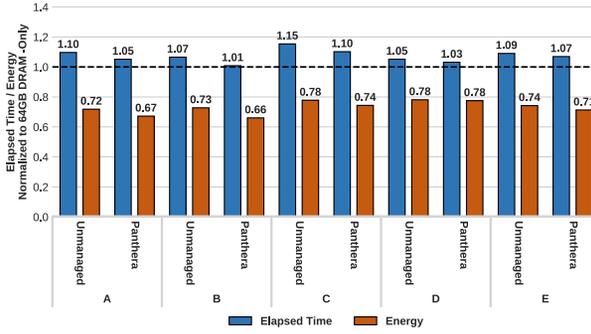


Fig. 16. Overall performance and energy results of QuickCached under a 64-GB heap; DRAM-to-memory ratio is 1/3.

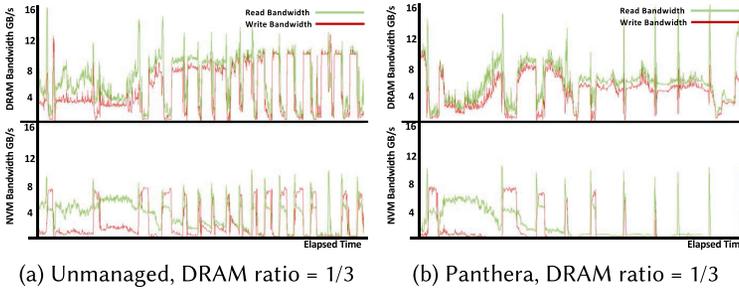


Fig. 17. GraphX-CC's memory access bandwidth.

7.4 Memory Access Analysis

NVM has high latency and low bandwidth. In general, the performance penalty caused by high latency increases with the number of memory accesses. For the same number of memory accesses, NVM incurs higher performance penalty for applications that have instantaneous bandwidth requirements which are beyond NVM's bandwidth. Figure 17 depicts the read/write bandwidth for unmanaged and Panthera on GraphX-CC. Compared to the unmanaged version, Panthera migrates most of the memory reads/writes from NVM to DRAM and reduces the high instantaneous memory access bandwidth requirements (i.e., peaks in the figure). Because Panthera allocates/moves frequently accessed data to DRAM, it reduces unnecessary NVM accesses (Sections 4.2.2 and 4.2.3).

7.5 Overhead of Monitoring and Migration

As discussed in Section 4.2, Panthera performs lightweight method-level monitoring on RDD objects to detect misplaced RDDs for dynamic migration. This subsection provides a closer examination of dynamic migration's overhead.

As we monitor only method calls invoked on RDD objects, we find dynamic monitoring overhead is negligible, i.e., it is less than 1% across our benchmarks. For example, for PageRank, only about 300 calls were observed on all RDD objects in a 20-minute execution. The second column of Table 5 reports the number of calls monitored for each application. For GraphX applications, which has thousands of RDD calls, the monitoring overheads are still less than 1%.

Dynamic migration (performed by the major GC) rarely occurs in our experiments, as can be seen from the third column of Table 5. There are two main reasons. First, the frequency of a major collection is very low because a majority of objects die young and most of the collection work

Table 5. Dynamic Monitoring and Migration

Program	# Calls monitored	# RDDs migrated
PR	328	0
KM	550	0
LR	333	0
TC	217	0
CC	2,945	1
SSSP	3,632	1
BC	336	0

is done by the minor GC. Second, for four applications (PR, KM, TC, and LR), our static analysis results are accurate enough and, hence, dynamic migration is never needed.

We observed that only two RDDs (during the executions of CC and SSSP) were migrated dynamically. Note that both CC and SSSP are GraphX applications. Each iteration of the processing creates new RDDs representing the updated graph and persists them. At the end of each iteration, the RDDs representing the old graph are explicitly *unpersisted*. Our static analysis, due to lack of support for the *unpersist* call, marks both old and new graph RDDs as hot data and generates a DRAM tag for all them. These RDD objects are then allocated in DRAM and their data objects are promoted eagerly to the DRAM space of the old generation. The RDD objects representing the old graphs, if they can survive a major GC, are migrated to the NVM space of the old generation due to their low access frequency.

To have better understanding of the individual contributions of pretenuring and dynamic migration, we have disabled the monitoring and migration and rerun the entire experiments. The performance difference was negligible (i.e., less than 1%). Hence, we conclude that most of Panthera’s benefit stems from pretenuring, which improves the performance of both the mutator and the GC. However, dynamic monitoring and migration increases the generality of Panthera’s optimizations, making Panthera applicable to applications with diverse access characteristics.

7.6 Performance and Energy of PGO

To examine the effectiveness of our profiling-guided optimization, we have evaluated Panthera with PGO enabled using the benchmarks listed in Table 4. We experimented with two heap sizes, 64 GB and 120 GB, and the nursery space is 1/6 of the heap size while 1/3 of the heap is DRAM. We compared Panthera w/ PGO against w/o PGO, and also against Unmanaged. Note our profiling is applied offline, thus it would not introduce extra overhead.

Figure 18 shows the overall performance and energy results when a 64-GB heap is used. The results are normalized w.r.t. those of the DRAM-only version. Compared with Panthera w/o PGO, the PGO can reduce the energy by 5.8% with 3.9% less execution time on average. Compared with 64-GB DRAM-only version, Panthera w/ PGO can achieve 36.4% energy reduction at 0.2% execution time overhead. However, we can see that some applications, such as LR, Graphx-CC, and Graphx-SSSP, can only get marginal benefits from PGO. There are two basic reasons. First, the access patterns in the RDDs of these applications are uniform and there is no need to divide the RDDs into finer chunks. In this situation, Panthera w/o PGO can recognize and migrate the data to correct place, as shown in Figure 7. Second, some memory accesses patterns, e.g., streaming, on the RDDs are easy to be recognized. The hardware and OS prefetching mechanisms work well for the data. In this case, even the Panthera w/ PGO can do a better data placement than Panthera w/o

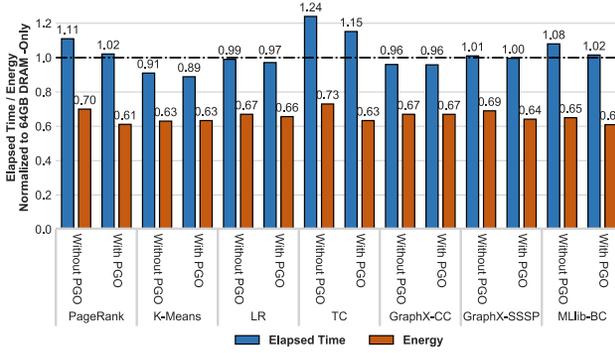


Fig. 18. Performance and energy results of Panthera with/without PGO. Heap size is 64 GB.

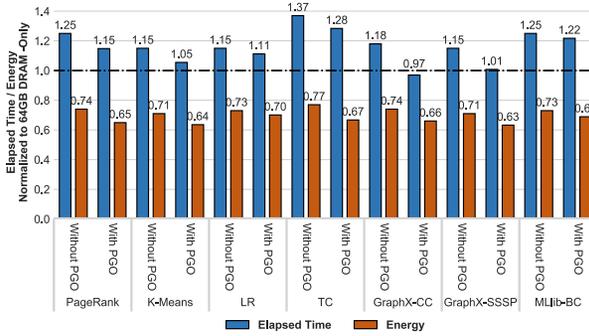


Fig. 19. Performance and energy results of Unmanaged with/without PGO. Heap size is 64 GB.

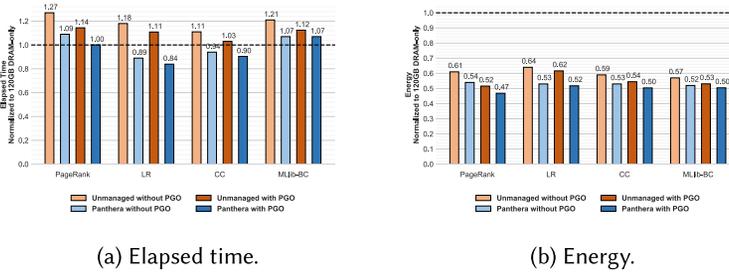


Fig. 20. Overall performance and energy results under a 120-GB heap.

PGO by recognizing and migrating chunks with higher access frequency to DRAM, it can't get too many benefits.

Our PGO can be decoupled from Panthera and be integrated into the unmanaged version, and Figure 19 shows the overall performance and energy results when PGO is implemented into unmanaged version. The results are normalized to the 64-GB DRAM-only version. Compared with Unmanaged w/o PGO, PGO reduces energy cost by 9.6% with 8.7% less execution time on average. Furthermore, compared with the 64-GB DRAM-only version, Unmanaged w/ PGO can achieve 33.8% energy reduction at 11.8% execution time overhead while Unmanaged w/o PGO reduces energy by 26.7% with a 21.4% execution time overhead.

Figure 20 shows the overall performance and energy results when a 120-GB heap is used. The results are normalized to DRAM-only version. With PGO, Unmanaged reduces energy by 44.8%

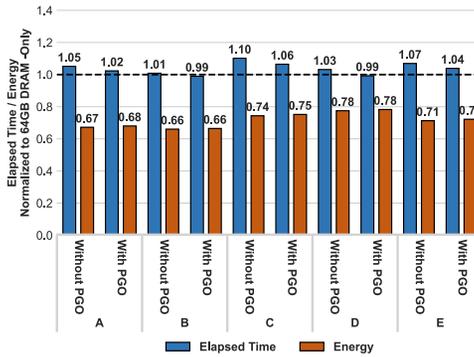


Fig. 21. Performance and energy results of Panthera with/without PGO.

with a 10.2% execution time overhead while Panthera reduces 50.1% energy with less than 1% execution time overhead on average.

PGO Results for QuickCached. Figure 21 reports the performance and energy results when employing PGO to QuickCached. We used the same experiment configuration as described in Section 7.2. The results are normalized w.r.t. those of the DRAM-only version. Compared with Panthera w/o PGO, the PGO reduced 2.9% execution time on average, while the energy consumption was almost the same. The benefits of PGO for QuickCached is less than for Spark, because: (1) the PGO was introduced to do finer-grained chunk placement and improve the performance further based on the static coarse-grained RDD-level analysis. For QuickCached, Panthera inferred the memory tag for each object at runtime, which is fine-grained already; and (2) the NVM access ratio for QuickCached is less than for Spark, e.g., accounting less than 10% of all the memory access, because of the biased object access characteristics as described in 2.5.

7.7 Discussion

How Far from the Ideal. For hybrid memories with managed runtime, the ideal solution is to place each object correctly according to its hotness, and adjust the placement dynamically considering the cost to do object migration and the benefit of migration. As Big Data systems usually have billions of objects in their heap during a normal execution, it is hard to profile the execution to record where each object should be placed, and when to do migration. Panthera introduced a simplified observation that we can develop a simple static analysis to infer the access pattern of each coarse-grained data collection, where all objects share the same pattern. In addition, this simple assumption is accurate enough as stated in 7.5.

The Effects of the Emulated NVM Specifications. Due to the limitations of Quartz, although most of the emulated latency/bandwidth results match the specifications used in the previous research [75], there is still some differences with the real NVM hardware, e.g., Optane DC [37]. There are two major differences between the emulated NVM and real NVM hardware. First, the emulated NVM doesn't show the asymmetry of NVM in read/write latency and bandwidth. Second, the emulated specifications is not exactly the same with the real hardware. For example, the read/write latencies of our emulated NVM are 300 ns/300 ns against the 305 ns/94 ns of the Optane DC [37]. However, we emphasize that different NVM technologies have different specifications and the emulated NVM already shows the performance characteristics of NVM and the performance difference between the DRAM and NVM. Our comprehensive evaluations on the emulated NVM can show the negligible overhead and high accuracy of our proposed static/dynamic analysis and

the effectiveness of our data migration policy. For example, we tuned the emulated NVM read/write latency from 300 ns/300 ns to 120 ns/120 ns, the baseline, Unmanaged, still has a significant performance degradation by suffering from the limited read/write bandwidth. Under these settings, Panthera still outperforms the baseline 11.6% on average. With a fixed 300-ns read/write latency, we also tuned the read/write bandwidth from 5 GB/s to 12 GB/s (12 GB/s is the bandwidth limit of our server QPI), Panthera always outperforms the baseline from 32.3% to 11.9% on average.

8 RELATED WORK

Hybrid Memories for Managed Runtime. To our knowledge, Panthera is the first practical work to optimize data layout in hybrid memories for managed-runtime-based distributed Big Data platforms. Existing efforts [11, 17, 32, 39, 40, 67, 72, 76, 80] that attempt to support persistent Java focus on regular applications or need to rebuild the platforms.

Inoue and Nakatani [36] identify code patterns in Java applications that can cause cache misses in L1 and L2. Gao et al. [31] propose a framework including support from hardware, the OS, and the runtime to extend NVM's lifetime. Two recent works close to Panthera are Espresso [80] and Write Rationing [11]. However, they were not designed for Big Data systems. Espresso is a JVM-based runtime system that enables persistent heaps. Developers can allocate objects in a persistent heap using a new instruction `pnew` while the runtime system provides crash consistency for the heap. Applying Espresso requires rewriting the Big Data platforms (e.g., Spark) using `pnew`, which is not practical.

Write Rationing [11] is a GC technique that places highly mutated objects in DRAM and mostly read objects in NVM to increase NVM lifetime. Like Espresso, this GC focuses on individual objects and does not consider application semantics. Panthera's nursery space is also placed in DRAM, similar to the Kingsguard-Nursery in Write Rationing. However, instead of focusing on individual objects, Panthera utilizes Spark semantics to obtain access information at the array granularity, leading to effective pretenuring and efficient runtime object tracking.

Memory Structure. There are two kinds of hybrid-memory structures: *flat structure*, where DRAM and NVM share a single memory space, and *vertical structure*, where DRAM is used as a buffer for NVM to store hot data. The vertical structure is normally managed by hardware and transparent to the OS and applications [44, 49, 54, 58, 69, 88, 90, 93]. Qureshi et al. [69] shows that a vertical structure with only 3% DRAM can reach similar performance to its DRAM-only version. However, the overhead of page monitoring and migration increases linearly with the working set [77]. The space overhead e.g., the tag store space of DRAM buffer, can also be high with a large volume of NVM [60].

Page-Based Migration. A great number of existing works use memory controllers to monitor page read/write frequency [20, 25, 30, 34, 53, 68, 70, 77, 88, 92] and migrate the top-ranked pages to DRAM. Another type of hybrid memory, composed of 3D-stacked DRAM and commodity DRAM, also adapts similar page monitoring policies [28, 38]. However, none of these techniques were designed for Big Data systems. Hassan et al. [34] show that, for some applications, migrating data at the object level can reduce power consumption.

For Big Data applications that have very large memory consumption, continuous monitoring at the page granularity can incur an unreasonable overhead. Page migration also incurs overhead in time and bandwidth. Bock et al. [18] report that page migration can increase execution time by 25% on average. Panthera uses static analysis to track memory usage at the RDD granularity, incorporating program semantics to reduce the dynamic monitoring overheads.

Static Data Placement. There exists a body of work that attempts to place data directly in appropriate spaces based either on their access frequencies [20, 53, 68, 74, 88] or on the result of a

program analysis [30, 34, 77]. Access frequency is normally calculated using a static data liveness analysis or offline profiling. Chatterjee et al. [20] place a single cache-line across multiple memory channels. Critical words (normally the first word) in a cache-line are placed in a low-latency channel. Wei et al. [77] show that the group of objects allocated by the same site in the source code exhibit similar lifetime behavior, which can be leveraged for static data placement. Dulloor et al. [30] classify memory accesses into three patterns and model the access time for a given mapping from the data structure with a specific access pattern to different memory types to get the optimal mapping configuration.

Li et al. [53] develop a binary instrumentation tool to statistically report memory access patterns in stack, heap, and global data. Phadke and Narayanasamy [68] profile an application's MLP and LLC misses to determine from which type of memory the application could benefit the most. Kim et al. [45] develop a key-value store for high-performance computers with large distributed NVM, which provides developers with a high-level interface to use the distributed NVM. However, none of these techniques were designed for managed Big Data systems.

9 CONCLUSION

We present Panthera, the first memory management technique for managed Big Data processing over hybrid memories. Panthera combines static analysis and GC techniques and profile-guided optimization to perform semantics-aware data placement in hybrid memory systems. Our evaluation shows that Panthera reduces energy significantly without incurring much extra time overhead.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their thorough and insightful comments. We are especially grateful to our shepherd Jennifer Sartor for her feedback, helping us improve the article substantially.

REFERENCES

- [1] 2012. LIBSVM Data: Classification. <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>.
- [2] 2017. Notre dame network dataset. <http://konect.uni-koblenz.de/networks/web-NotreDame>.
- [3] 2017. QuickCached. <https://github.com/QuickServerLab/QuickCached>.
- [4] 2017. Wikipedia links, network dataset. <http://konect.uni-koblenz.de/networks>.
- [5] 2018. 3D XPointTM: A Breakthrough in Non-Volatile Memory Technology. <https://www.intel.com/content/www/us/en/architecture-and-technology/intel-micron-3d-xpoint-webcast.html>.
- [6] 2019. Apache SparkTM. <https://spark.apache.org>.
- [7] 2019. Intel VTuneTM Amplifier. <https://software.intel.com/en-us/vtune>.
- [8] 2019. OpenJDK. <https://openjdk.java.net>.
- [9] Neha Agarwal and Thomas F. Wenisch. 2017. Thermostat: Application-transparent page management for two-tiered main memory. In *Proceedings of the 22nd International Conference on Architectural Support for Programming Languages and Operating Systems*. 631–644.
- [10] Shoaib Akram, Jennifer B. Sartor, Kathryn S. McKinley, and Lieven Eeckhout. 2018. Emulating hybrid memory on NUMA hardware. *CoRR* (2018).
- [11] Shoaib Akram, Jennifer B. Sartor, Kathryn S. McKinley, and Lieven Eeckhout. 2018. Write-rationing garbage collection for hybrid memories. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI'18)*. ACM, New York, 62–77.
- [12] Emmanuel Amaro, Christopher Branner-Augmon, Zhihong Luo, Amy Ousterhout, Marcos K. Aguilera, Aurojit Panda, Sylvia Ratnasamy, and Scott Shenker. 2020. Can far memory improve job throughput? In *Proceedings of the 15th European Conference on Computer Systems (EuroSys'20)*. ACM, New York, Article 14, 16 pages. <https://doi.org/10.1145/3342195.3387522>
- [13] Joy Arulraj, Justin Levandoski, Umar Farooq Minhas, and Per-Ake Larson. 2018. Bztree: A high-performance latch-free range index for non-volatile memory. *Proc. VLDB Endow.* 11, 5 (Jan. 2018), 553–565.
- [14] Joy Arulraj, Andrew Pavlo, and Subramanya R. Dulloor. 2015. Let's talk about storage & recovery methods for non-volatile memory database systems. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. ACM, New York, 707–722.

- [15] Joy Arulraj, Matthew Perron, and Andrew Pavlo. 2016. Write-behind logging. *Proc. VLDB Endow.* 10, 4 (Nov. 2016), 337–348.
- [16] Berk Atikoglu, Yuehai Xu, Eitan Frachtenberg, Song Jiang, and Mike Paleczny. 2012. Workload analysis of a large-scale key-value store. In *Proceedings of the 12th ACM SIGMETRICS/PERFORMANCE Joint International Conference on Measurement and Modeling of Computer Systems*. ACM, New York, 53–64.
- [17] M. P. Atkinson, L. Daynès, M. J. Jordan, T. Printezis, and S. Spence. 1996. An orthogonally persistent Java. *SIGMOD Rec.* 25, 4 (Dec. 1996), 68–75.
- [18] Santiago Bock, Bruce R. Childers, Rami G. Melhem, and Daniel Mossé. 2014. Concurrent page migration for mobile systems with OS-managed hybrid memory. In *Proceedings of the 11th ACM Conference on Computing Frontiers (CF'14)*. ACM, New York, 31:1–31:10.
- [19] Sergey Brin and Lawrence Page. 1998. The anatomy of a large-scale hypertextual web search engine. *Comput. Netw. ISDN Syst.* 30, 1–7 (April 1998), 107–117.
- [20] N. Chatterjee, M. Shevgoor, R. Balasubramonian, A. Davis, Z. Fang, R. Illikkal, and R. Iyer. 2012. Leveraging heterogeneity in DRAM main memories to accelerate critical word access. In *Proceedings of the 45th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-45)*. ACM, New York, 13–24.
- [21] Youmin Chen, Youyou Lu, Fan Yang, Qing Wang, Yang Wang, and Jiwu Shu. 2020. FlatStore: An efficient log-structured key-value storage engine for persistent memory. In *Proceedings of the 25th International Conference on Architectural Support for Programming Languages and Operating Systems*. 1077–1091.
- [22] Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. 2010. Benchmarking cloud serving systems with YCSB. In *Proceedings of the 1st ACM Symposium on Cloud Computing*. ACM, New York, 143–154.
- [23] Intel Corporation. 2015. An introduction to pmemcheck. <https://pmem.io/2015/07/17/pmemcheck-basic.html>.
- [24] Intel Corporation. 2018. Redis. <https://github.com/pmem/redis/tree/3.2-nvml>.
- [25] G. Dhiman, R. Ayoub, and T. Rosing. 2009. PDRAM: A hybrid PRAM and DRAM main memory system. In *Proceedings of the 46th Annual Design Automation Conference (DAC'09)*. ACM, New York, 664–669.
- [26] Bang Di, Jiawen Liu, Hao Chen, and Dong Li. 2021. Fast, flexible, and comprehensive bug detection for persistent memory programs. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*. ACM, New York, 503–516.
- [27] Mingkai Dong, Heng Bu, Jifei Yi, Benchao Dong, and Haibo Chen. 2019. Performance and protection in the ZoFS user-space NVM file system. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*. ACM, New York, 478–493.
- [28] Xiangyu Dong, Yuan Xie, Naveen Muralimanohar, and Norman P. Jouppi. 2010. Simple but effective heterogeneous main memory with on-chip memory controller support. In *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis (SC'10)*. ACM, New York, 1–11.
- [29] Subramanya R. Dulloor, Sanjay Kumar, Anil Keshavamurthy, Philip Lantz, Dheeraj Reddy, Rajesh Sankaran, and Jeff Jackson. 2014. System software for persistent memory. In *Proceedings of the 9th European Conference on Computer Systems (EuroSys'14)*. 15:1–15:15.
- [30] 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 11th European Conference on Computer Systems (EuroSys'16)*. 15:1–15:16.
- [31] Tiejun Gao, Karin Strauss, Stephen M. Blackburn, Kathryn S. McKinley, Doug Burger, and James R. Larus. 2013. Using managed runtime systems to tolerate holes in wearable memories. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI'13)*. ACM, New York, 297–308.
- [32] Lokesh Gidra, Gaël Thomas, Julien Sopena, Marc Shapiro, and Nhan Nguyen. 2015. NumaGiC: A garbage collector for big data on big NUMA machines. In *Proceedings of the 20th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'15)*. 661–673.
- [33] Joseph E. Gonzalez, Reynold S. Xin, Ankur Dave, Daniel Crankshaw, Michael J. Franklin, and Ion Stoica. 2014. GraphX: Graph processing in a distributed dataflow framework. In *Proceedings of the 11th USENIX Conference on Operating Systems Design and Implementation (OSDI'14)*. 599–613.
- [34] Ahmad Hassan, Hans Vandierendonck, and Dimitrios S. Nikolopoulos. 2015. Software-managed energy-efficient hybrid DRAM/NVM main memory. In *Proceedings of the 12th ACM International Conference on Computing Frontiers (CF'15)*. ACM, New York, 23:1–23:8.
- [35] Mark Hildebrand, Jawad Khan, Sanjeev Trika, Jason Lowe-Power, and Venkatesh Akella. 2020. AUTOTM: Automatic tensor movement in heterogeneous memory systems using integer linear programming. In *Proceedings of the 25th International Conference on Architectural Support for Programming Languages and Operating Systems*. 875–890.
- [36] Hiroshi Inoue and Toshio Nakatani. 2012. Identifying the sources of cache misses in Java programs without relying on hardware counters. In *Proceedings of the 2012 International Symposium on Memory Management (ISMM'12)*. 133–142.

- [37] Joseph Izraelevitz, Jian Yang, Lu Zhang, Juno Kim, Xiao Liu, Amirsaman Memaripour, Yun Joon Soh, Zixuan Wang, Yi Xu, Subramanya R. Dulloor, Jishen Zhao, and Steven Swanson. 2019. Basic Performance Measurements of the Intel Optane DC Persistent Memory Module. [arXiv:cs/1903.05714](https://arxiv.org/abs/cs/1903.05714).
- [38] X. Jiang, N. Madan, L. Zhao, M. Upton, R. Iyer, S. Makineni, D. Newell, Y. Solihin, and R. Balasubramonian. 2010. CHOP: Adaptive filter-based DRAM caching for CMP server platforms. In *Proceedings of the 16th International Symposium on High-Performance Computer Architecture (HPCA'10)*. 1–12.
- [39] Mick Jordan. 1996. Early experiences with persistent Java. In *The First International Workshop on Persistence and Java*.
- [40] Mick Jordan and Malcolm Atkinson. 2000. *Orthogonal Persistence for the JavaTM Platform: Specification and Rationale*. Technical Report. Mountain View, CA.
- [41] Rohan Kadekodi, Se Kwon Lee, Sanidhya Kashyap, Taesoo Kim, Aasheesh Kolli, and Vijay Chidambaram. 2019. SplitFS: Reducing software overhead in file systems for persistent memory. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*. ACM, New York, 494–508.
- [42] Olzhas Kaiyrakhmet, Songyi Lee, Beomseok Nam, Sam H. Noh, and Young-ri Choi. 2019. SLM-DB: Single-level key-value store with persistent memory. In *Proceedings of the 17th {USENIX} Conference on File and Storage Technologies ({FAST} 19)*. 191–205.
- [43] Sudarsun Kannan, Ada Gavrilovska, Vishal Gupta, and Karsten Schwan. 2017. HeteroOS: OS design for heterogeneous memory management in datacenter. In *Proceedings of the 44th Annual International Symposium on Computer Architecture*. 521–534.
- [44] Taeho Kgil, David Roberts, and Trevor Mudge. 2008. Improving NAND flash based disk caches. In *Proceedings of the 35th Annual International Symposium on Computer Architecture (ISCA'08)*. 327–338.
- [45] Jungwon Kim, Seyong Lee, and Jeffrey S. Vetter. 2017. PapyrusKV: A high-performance parallel key-value store for distributed NVM architectures. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC'17)*. 57:1–57:14.
- [46] Emre Kultursay, Mahmut Kandemir, Anand Sivasubramaniam, and Onur Mutlu. 2013. Evaluating STT-RAM as an energy-efficient main memory alternative. In *Proceedings of the IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS'13)*. 256–267.
- [47] Harendra Kumar, Yuvraj Patel, Ram Kesavan, and Sumith Makam. 2017. High performance metadata integrity protection in the {WAFL} copy-on-write file system. In *Proceedings of the 15th {USENIX} Conference on File and Storage Technologies ({FAST} 17)*. 197–212.
- [48] Youngjin Kwon, Henrique Fingler, Tyler Hunt, Simon Peter, Emmett Witchel, and Thomas Anderson. 2017. Strata: A cross media file system. In *Proceedings of the 26th Symposium on Operating Systems Principles*. 460–477.
- [49] Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger. 2009. Architecting phase change memory as a scalable DRAM alternative. In *Proceedings of the 36th Annual International Symposium on Computer Architecture (ISCA'09)*. 2–13.
- [50] Benjamin C. Lee, Ping Zhou, Jun Yang, Youtao Zhang, Bo Zhao, Engin Ipek, Onur Mutlu, and Doug Burger. 2010. Phase-change technology and the future of main memory. *IEEE Micro* 30, 1 (Jan. 2010), 143–143.
- [51] Se Kwon Lee, Jayashree Mohan, Sanidhya Kashyap, Taesoo Kim, and Vijay Chidambaram. 2019. Recipe: Converting concurrent DRAM indexes to persistent-memory indexes. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*. 462–477.
- [52] Lenovo. 2018. Memcached-pmem. <https://github.com/lenovo/memcachedpmem>.
- [53] Dong Li, Jeffrey S. Vetter, Gabriel Marin, Collin McCurdy, Cristian Cira, Zhuo Liu, and Weikuan Yu. 2012. Identifying opportunities for byte-addressable non-volatile memory in extreme-scale scientific applications. In *Proceedings of the 2012 IEEE 26th International Parallel and Distributed Processing Symposium (IPDPS'12)*. 945–956.
- [54] Yang Li, Saugata Ghose, Jongmoo Choi, Jin Sun, Hui Wang, and Onur Mutlu. 2017. Utility-based hybrid memory management. In *Proceedings of the 2017 IEEE International Conference on Cluster Computing (CLUSTER'17)*. 152–165.
- [55] Sihang Liu, Korakit Seemakhupt, Yizhou Wei, Thomas Wenisch, Aasheesh Kolli, and Samira Khan. 2020. Cross-failure bug detection in persistent memory programs. In *Proceedings of the 25th International Conference on Architectural Support for Programming Languages and Operating Systems*. 1187–1202.
- [56] Baotong Lu, Xiangpeng Hao, Tianzheng Wang, and Eric Lo. 2020. Dash: Scalable hashing on persistent memory. *arXiv preprint arXiv:2003.07302* (2020).
- [57] Martin Maas, Krste Asanović, Tim Harris, and John Kubiawicz. 2016. Taurus: A holistic language runtime system for coordinating distributed managed-language applications. In *Proceedings of the 21st International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'16)*. 457–471.
- [58] Prasanth Mangalagiri, Karthik Sarpatwari, Aditya Yanamandra, VijayKrishnan Narayanan, Yuan Xie, Mary Jane Irwin, and Osama Awadel Karim. 2008. A low-power phase change memory based hybrid cache architecture. In *Proceedings of the 18th ACM Great Lakes Symposium on VLSI (GLSVLSI'08)*. 395–398.

- [59] Hasan Al Maruf and Mosharaf Chowdhury. 2020. Effectively prefetching remote memory with leap. In *Proceedings of the 2020 USENIX Annual Technical Conference (USENIX ATC 20)*. USENIX Association, 843–857. <https://www.usenix.org/conference/atc20/presentation/al-maruf>.
- [60] Justin Meza, Jichuan Chang, HanBin Yoon, Onur Mutlu, and Parthasarathy Ranganathan. 2012. Enabling efficient and scalable hybrid memories using fine-granularity DRAM cache management. *IEEE Computer Architecture Letters* 11, 2 (July 2012), 61–64.
- [61] Micron. 2017. TN-40-07: Calculating Memory Power for DDR4 SDRAM Introduction. https://www.micron.com/-/media/documents/products/technical-note/dram/tn4007_ddr4_power_calculation.pdf.
- [62] Jeffrey C. Mogul, Eduardo Argollo, Mehul Shah, and Paolo Faraboschi. 2009. Operating system support for NVM+DRAM hybrid main memory. In *Proceedings of the 12th Conference on Hot Topics in Operating Systems (HotOS'09)*. 14–14.
- [63] Gaku Nakagawa and Shuichi Oikawa. 2015. NVM/DRAM hybrid memory management with language runtime support via MRW queue. In *Proceedings of the 2015 IEEE/ACIS 16th International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD'15)*. 357–362.
- [64] Moohyeon Nam, Hokeun Cha, Young-ri Choi, Sam H. Noh, and Beomseok Nam. 2019. Write-optimized dynamic hashing for persistent memory. In *Proceedings of the 17th {USENIX} Conference on File and Storage Technologies ({FAST} 19)*. 31–44.
- [65] Khanh Nguyen, Lu Fang, Guoqing Xu, Brian Demsky, Shan Lu, Sanazsadat Alamian, and Onur Mutlu. 2016. Yak: A high-performance big-data-friendly garbage collector. In *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation (OSDI'16)*. 349–365.
- [66] Khanh Nguyen, Kai Wang, Yingyi Bu, Lu Fang, Jianfei Hu, and Guoqing Xu. 2015. FACADE: A compiler and runtime for (almost) object-bounded big data applications. In *Proceedings of the 20th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'15)*. 675–690.
- [67] James O'Toole, Scott Nettles, and David Gifford. 1993. Concurrent compacting garbage collection of a persistent heap. In *Proceedings of the 14th ACM Symposium on Operating Systems Principles (SOSP'93)*. ACM, New York, 161–174.
- [68] Sujay Phadke and Satish Narayanasamy. 2011. MLP aware heterogeneous memory system. In *Proceedings of 2011 IEEE Design, Automation Test Conference in Europe (DATE'11)*. 1–6.
- [69] Moinuddin K. Qureshi, Vijayalakshmi Srinivasan, and Jude A. Rivers. 2009. Scalable high performance main memory system using phase-change memory technology. In *Proceedings of the 36th Annual International Symposium on Computer Architecture (ISCA'09)*. 24–33.
- [70] Luiz E. Ramos, Eugene Gorbato, and Ricardo Bianchini. 2011. Page placement in hybrid memory systems. In *Proceedings of the International Conference on Supercomputing (ICS'11)*. 85–95.
- [71] Jinglei Ren, Jishen Zhao, Samira Khan, Jongmoo Choi, Yongwei Wu, and Onur Mutlu. 2015. ThyNVM: Enabling software-transparent crash consistency in persistent memory systems. In *Proceedings of the 2015 48th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 672–685.
- [72] M. Satyanarayanan, Henry H. Mashburn, Puneet Kumar, David C. Steere, and James J. Kistler. 1994. Lightweight recoverable virtual memory. *ACM Trans. Comput. Syst.* 12, 1 (Feb. 1994), 33–57.
- [73] Alexander van Renen, Viktor Leis, Alfons Kemper, Thomas Neumann, Takushi Hashida, Kazuichi Oe, Yoshiyasu Doi, Lilian Harada, and Mitsuru Sato. 2018. Managing non-volatile memory in database systems. In *Proceedings of the 2018 International Conference on Management of Data*. 1541–1555.
- [74] Nandita Vijaykumar, Abhilasha Jain, Diptesh Majumdar, Kevin Hsieh, Gennady Pekhimenko, Eiman Ebrahimi, Nas-taran Hajinazar, Phillip B. Gibbons, and Onur Mutlu. 2018. A case for richer cross-layer abstractions: Bridging the semantic gap with expressive memory. In *Proceedings of the 45th Annual International Symposium on Computer Architecture (ISCA'18)*. 207–220.
- [75] Haris Volos, Guilherme Magalhaes, Ludmila Cherkasova, and Jun Li. 2015. Quartz: A lightweight performance emulator for persistent memory software. In *Proceedings of the 16th Annual Middleware Conference (Middleware'15)*. 37–49.
- [76] Chenxi Wang, Ting Cao, John Zigman, Fang Lv, Yunquan Zhang, and Xiaobing Feng. 2016. Efficient management for hybrid memory in managed language runtime. In *Proceedings of the 16th IFIP International Conference on Network and Parallel Computing (NPC'16)*. 29–42.
- [77] Wei Wei, Dejun Jiang, Sally A. McKee, Jin Xiong, and Mingyu Chen. 2015. Exploiting program semantics to place data in hybrid memory. In *Proceedings of the 2015 International Conference on Parallel Architecture and Compilation (PACT'15)*. 163–173.
- [78] H. S. P. Wong, H. Lee, S. Yu, Y. Chen, Y. Wu, P. Chen, B. Lee, F. T. Chen, and M. Tsai. 2012. Metal-oxide RRAM. *Proc. IEEE* 100, 6 (June 2012), 1951–1970.
- [79] H. S. P. Wong, S. Raoux, S. Kim, J. Liang, J. P. Reifenberg, B. Rajendran, M. Asheghi, and K. E. Goodson. 2010. Phase change memory. *Proc. IEEE* 98, 12 (Dec 2010), 2201–2227.

- [80] Mingyu Wu, Ziming Zhao, Haoyu Li, Heting Li, Haibo Chen, Binyu Zang, and Haibing Guan. 2018. Espresso: Brewing Java for more non-volatility. In *Proceedings of the 20th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'18)*. 70–83.
- [81] Xingbo Wu, Fan Ni, Li Zhang, Yandong Wang, Yufei Ren, Michel Hack, Zili Shao, and Song Jiang. 2016. NVMcached: An NVM-based key-value cache. In *Proceedings of the 7th ACM SIGOPS Asia-Pacific Workshop on Systems*. 1–7.
- [82] Fei Xia, Dejun Jiang, Jin Xiong, and Ninghui Sun. 2017. HiKV: A hybrid index key-value store for DRAM-NVM memory systems. In *Proceedings of the 2017 {USENIX} Annual Technical Conference ({USENIX}{ATC} 17)*. 349–362.
- [83] Jian Xu and Steven Swanson. 2016. {NOVA}: A log-structured file system for hybrid volatile/non-volatile main memories. In *Proceedings of the 14th {USENIX} Conference on File and Storage Technologies ({FAST} 16)*. 323–338.
- [84] Jian Xu, Lu Zhang, Amirsaman Memaripour, Akshatha Gangadharaiah, Amit Borase, Tamires Brito Da Silva, Steven Swanson, and Andy Rudoff. 2017. Nova-fortis: A fault-tolerant non-volatile main memory file system. In *Proceedings of the 26th Symposium on Operating Systems Principles*. 478–496.
- [85] Zi Yan, Daniel Lustig, David Nellans, and Abhishek Bhattacharjee. 2019. Nimble page management for tiered memory systems. ACM, New York. <https://doi.org/10.1145/3297858.3304024>
- [86] Zi Yan, Daniel Lustig, David Nellans, and Abhishek Bhattacharjee. 2019. Translation ranger: Operating system support for contiguity-aware TLBs. In *Proceedings of the 46th International Symposium on Computer Architecture*. 698–710.
- [87] Yanfei Yang, Mingyu Wu, Haibo Chen, and Binyu Zang. 2021. Bridging the performance gap for copy-based garbage collectors atop non-volatile memory. ACM, New York. <https://doi.org/10.1145/3447786.3456246>
- [88] HanBin Yoon, Justin Meza, Rachata Ausavarungnirun, Rachael Harding, and Onur Mutlu. 2012. Row buffer locality aware caching policies for hybrid memories. In *Proceedings of the 2012 IEEE 30th International Conference on Computer Design (ICCD'12)*. 337–344.
- [89] Hanbin Yoon, Justin Meza, Naveen Muralimanohar, Norman P. Jouppi, and Onur Mutlu. 2014. Efficient data mapping and buffering techniques for multilevel cell phase-change memories. *ACM Trans. Archit. Code Optim.* 11, 4 (Dec. 2014), 40:1–40:25.
- [90] Xiangyao Yu, Christopher J. Hughes, Nadathur Satish, Onur Mutlu, and Srinivas Devadas. 2017. Banshee: Bandwidth-efficient DRAM caching via software/hardware cooperation. In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-50)*. ACM, New York, 1–14.
- [91] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauly, Michael J. Franklin, Scott Shenker, and Ion Stoica. 2012. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *Presented as part of the 9th USENIX Symposium on Networked Systems Design and Implementation (NSDI'12)*. 15–28.
- [92] Wangyuan Zhang and Tao Li. 2009. Exploring phase change memory and 3d die-stacking for power/thermal friendly, fast and durable memory architectures. In *Proceedings of the 2009 18th International Conference on Parallel Architectures and Compilation Techniques (PACT'09)*. 101–112.
- [93] Ping Zhou, Bo Zhao, Jun Yang, and Youtao Zhang. 2009. A durable and energy efficient main memory using phase change memory technology. In *Proceedings of the 36th Annual International Symposium on Computer Architecture (ISCA'09)*. 14–23.
- [94] Omer Zilberberg, Shlomo Weiss, and Sivan Toledo. 2013. Phase-change memory: An architectural perspective. *ACM Comput. Surv.* 45, 3 (July 2013), 29:1–29:33.
- [95] Pengfei Zuo, Yu Hua, and Jie Wu. 2018. Write-optimized and high-performance hashing index scheme for persistent memory. In *Proceedings of the 13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*. 461–476.
- [96] Yoav Zuriel, Michal Friedman, Gali Sheffi, Nachshon Cohen, and Erez Petrank. 2019. Efficient lock-free durable sets. In *Proceedings of the ACM on Programming Languages 3, (OOPSLA 2019)*. ACM, New York, 1–26.

Received November 2020; revised October 2021; accepted January 2022