

Canvas: Isolated and Adaptive Swapping for Multi-Applications on Remote Memory

Chenxi Wang^{†*} Yifan Qiao^{†*} Haoran Ma[†] Shi Liu[†] Yiyang Zhang[‡]
Wenguang Chen[§] Ravi Netravali[‡] Miryung Kim[†] Guoqing Harry Xu[†]
UCLA[†] UCSD[‡] Tsinghua University[§] Princeton University[‡]

Abstract

Remote memory techniques for datacenter applications have recently gained a great deal of popularity. Existing remote memory techniques focus on the efficiency of a single application setting only. However, when multiple applications co-run on a remote-memory system, significant interference could occur, resulting in unexpected slowdowns even if the same amounts of physical resources are granted to each application. This slowdown stems from massive sharing in applications’ swap data paths. Canvas is a redesigned swap system that fully isolates swap paths for remote-memory applications. Canvas allows each application to possess its dedicated swap partition, swap cache, prefetcher, and RDMA bandwidth. Swap isolation lays a foundation for adaptive optimization techniques based on each application’s own access patterns and needs. We develop three such techniques: (1) adaptive swap entry allocation, (2) semantics-aware prefetching, and (3) two-dimensional RDMA scheduling. A thorough evaluation with a set of widely-deployed applications demonstrates that Canvas minimizes performance variation and dramatically reduces performance degradation.

1 Introduction

Techniques enabling datacenter applications to use far memory [34, 37, 8, 59, 70, 87, 100, 86, 17] have gained traction due to their potential to break servers’ memory capacity wall, thereby improving performance and resource utilization. Existing far-memory techniques can be roughly classified into two categories: (1) clean-slate techniques [86, 17] that provide new primitives for developers to manage remote memory, and (2) swap-based techniques [37, 87, 8, 100, 2] that piggyback on existing swap mechanisms in the OS kernel. Clean-slate techniques provide greater efficiency by enabling user-space far memory accesses, while swap-based techniques offer transparency, allowing legacy code to run *as is* on a far-memory system. This paper focuses on swap mechanisms as they are more practical and easier to adopt.

A typical swap system in the OS uses a *swap partition* and *swap cache* for applications to swap data between memory and external storage. The swap partition is a storage-backed swap space. The swap cache is an intermediate buffer between the *local memory* and storage—it caches *unmapped*

pages that were just swapped in or are about to be swapped out. Upon a page fault, the OS looks up the swap cache; a cache miss would trigger a *demand swap* and a number of *prefetching swaps*. Swaps are served by RDMA and all fetched pages are initially placed in the swap cache. The demand page is then mapped to a virtual page and moved out of the swap cache, completing the fault handling process.

Problems. Current swap systems run multiple applications over shared swap resources (*i.e.*, swap partition, RDMA, *etc.*). This design works for *disk-based swapping* where disk access is slow—each application can allow only a tiny number of pages to be swapped to maintain an acceptable overhead. This assumption, however, no longer holds under far memory because an application can place more data in far memory than local memory and yet still be efficient, thanks to RDMA’s low latency and high bandwidth.

As such, applications have orders-of-magnitude more swap requests under far memory than disks. Millions of swap requests from different applications go through the same shared data path in a short period of time, leading to *severe performance interference*. Our experiments show that, with the same amounts of CPU and local-memory resources, co-running applications leads up to a 6.4× slowdown, an overhead unacceptable for any real-world deployment.

State of the Art. Interference is a known problem in datacenter applications and a large body of work exists on isolation of CPU [61, 14, 23], I/O [38, 92], network bandwidth [12, 35, 90, 83, 74, 50] and processing [56]. Most of these techniques build on Linux’s *cgroup* mechanism, which focuses on isolation of traditional resources such as CPU and memory, *not* swap resources such as remote memory usage and RDMA. Prior swap optimizations such as *Infiniswap* [37], *Fastswap* [8], or *NVMe-over-fabrics* [2] focus on reducing remote access latency, overlooking the impact of swap interference in realistic settings.

Contribution #1: Interference Study (§3). We conducted a systematic study with a set of widely-deployed applications on Linux 5.5, the latest kernel version compatible with Mellanox’s latest driver (4.9-3.1.5.0) for our InfiniBand card. Our results reveal three major performance problems:

- **Severe lock contention:** Since all applications share a single swap partition, extensive locking is needed for swap entry allocation (needed by every swap-out), reducing throughput and precluding full utilization of RDMA’s bandwidth. Our experience shows that in windows of fre-

* Authors contributed equally

quent remote accesses, applications can spend **70%** of the windows’ time on swap entry allocation.

- **Uncontrolled use of swap resources (e.g., RDMA):** The use of the shared RDMA bandwidth is often dominated by the pages fetched for applications with many threads simultaneously performing frequent remote accesses. For example, aggressively (pre)fetching pages to fulfill one application’s needs can disproportionately reduce other applications’ bandwidth usage. Further, even within one application, prefetching competes resources with demand swaps, leading to either prolonged fault handling or delayed prefetching that fails to bring back pages in time.
- **Reduced prefetching accuracy:** Applications use the same prefetcher, prefetching data based on *low-level (sequential or strided) access patterns* across applications. However, modern applications exhibit far more diverse access patterns, making it hard for prefetching to be effective across the board. For example, co-running Spark and native applications reduces LEAP [70]’s prefetching contribution by **3.19×**.

These results highlight two main problems. First, interference is caused by sharing a combination of swap resources including the swap partition/cache, and RDMA (bandwidth and SRAM on RNIC). Although recent kernel versions added support [45] for charging prefetched pages into `cgroup`, resolving interference requires a *holistic* approach that can isolate all these resources. Furthermore, interference stems not only from resource racing, but also from fundamental limitations with the current design of the swap system. For instance, reducing interference between prefetching and demand swapping requires understanding whether a prefetching request can come back in time. If not, it should be dropped to give resources to demand requests, which are on the critical path. This, in turn, requires a re-design of the kernel’s fault handling logic.

Second, cloud applications exhibit highly diverse behaviors and resource profiles. For example, applications with a great number of threads are more sensitive to locking than single-threaded applications. Furthermore, managed applications such as Spark often make heavy use of reference-based data structures while native applications are often dominated by large arrays. The *application-agnostic nature* of the swap system makes it hard for a one-size-fits-all policy (e.g., a global prefetcher) to work well for diverse applications. Effective per-application policies dictates (1) holistic swap isolation and (2) understanding application semantics, which is currently inaccessible in the kernel.

Contribution #2: Holistic Swap Isolation (§4). To solve the first problem, we develop Canvas, a *fully-isolated* swap system, which enables each application to have its dedicated swap partition, swap cache, and RDMA usage. In doing so, Canvas can charge each application’s `cgroup` for the usage

of all kinds of swap resources, preventing certain applications from aggressively invading others’ resources.

Contribution #3: Isolation-Enabled Adaptive Optimizations (§5). To solve the second problem, we develop a set of adaptive optimizations that can tailor their policies and strategies to application-specific swap behaviors and resource needs. Our adaptive optimizations bring a *further boost* on top of the isolation-provided benefits, making co-running applications even *outperform* their individual runs.

(1) Adaptive Swap Entry Allocation (§5.1) Separating swap partitions reduces lock contention at swap entry allocations to a certain degree, but the contention can still be heavy for multi-threaded applications. For example, Spark creates many threads to fully utilize cores and these threads need synchronizations before obtaining swap entries. The synchronization overhead increases dramatically with the number of cores (§6.3.1), creating a scalability bottleneck. We develop an adaptive swap entry allocator that dynamically balances between the degree of lock contention (i.e., time) and the amount of swap space needed (i.e., space) based on each application’s memory behaviors.

(2) Adaptive Two-tier Prefetching (§5.2) Current kernel prefetchers build on low-level access patterns (e.g., sequential or strided). Although such patterns are useful for applications with large array usages, many cloud applications are written in high-level, managed languages such as Java or Python; their accesses come from multiple threads or exhibit pointer-chasing behavior as opposed to sequential or strided patterns. As effective prefetching is paramount to remote-memory performance, Canvas employs a two-tier prefetching design. Our *kernel-tier prefetcher* prefetches data for each application into its private swap cache based on low-level patterns. Once this prefetcher cannot effectively prefetch data, Canvas adaptively forwards the faulty address up to the *application tier* via a modified `userfaultfd` interface, enabling customized prefetching logic at the level of reference-based or thread-based access patterns.

(3) Adaptive RDMA Scheduling (§5.3) Isolating RDMA bandwidth alone for each application is *not* sufficient. Because there could be many more *prefetching* requests than *demand swap requests*, naively sending all to RDMA delays demand requests, increasing fault-handling latency. On the other hand, naively delaying prefetching requests (as in FastSwap [8]) reduces their *timeliness*, making prefetched pages useless. We build a *two-dimensional* RDMA scheduler, which schedules packets not only between applications but also between prefetching and demand requests for each application.

Results. We implemented Canvas in Linux 5.5. We also modified Oracle’s OpenJDK 12 to add support for application-tier prefetching for high-level languages. An evaluation (§6) with a set of 12 widely-deployed applications (including Apache Spark [104], Memcached [4], XGBoost [21, 20], Snappy [36], etc.) demonstrates that Canvas

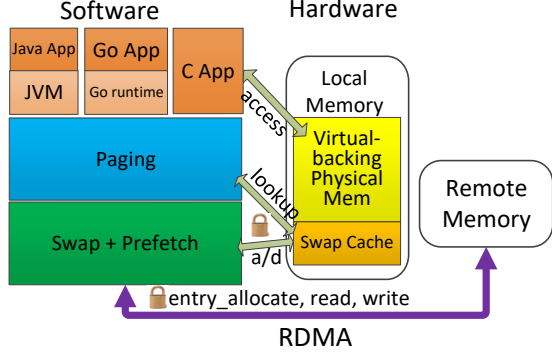


Figure 1: Data plane and remote-access data path; a majority of local memory backs virtual memory (yellow), while a small portion of it is used by the swap cache (orange).

improves the overall application performance by up to $6.12\times$ (average $1.84\times$) and reduces applications’ performance variation (*i.e.*, standard deviation) by $31\times$, from an overall of 1.53 to 0.05 . Canvas increases the overall RDMA bandwidth by $2.8\times$ for co-run applications, and outperforms the allocator in the latest kernel version (5.14) by $13\times$ when Memcached runs on RAMDisk with 48 cores.

2 Background

This section presents the necessary background in Linux 5.5, which is the latest kernel version compatible with Mellanox’s latest driver for our InfiniBand card.

Figure 1 illustrates the kernel’s data plane for remote-memory applications where remote memory is mapped into the host server as a swap partition. The swap partition is split into a set of 4KB *swap entries*, a unit of management (*e.g.*, creation, read, and write). Each swap entry maps to an actual remote memory cell and has a unique entry ID. Applications access the local memory. Upon a page fault, the kernel uses the swap entry ID contained in the corresponding page table entry (PTE) to locate the swap entry that stores the page.

Before issuing a remote fetch, the kernel looks up the swap cache, which is a set of radix trees, each containing a number of cached and unmapped pages for a block (*e.g.*, 64MB) of swap entries. These pages were either just swapped in due to demand swapping or prefetching, or are about to be swapped out. If a page can be found there, it gets mapped to a virtual page and removed from the swap cache. Otherwise, the kernel issues an I/O request, which is then pushed into RDMA’s dispatch queue. When a demand swap occurs, the kernel prefetches a number of pages that will likely be needed in the future. This number depends on the swap history at the past few page faults. For example, if the pages fetched follow a sequential or strided pattern, the kernel will follow this pattern to fetch a few more pages. If no pattern is found, the kernel reduces the number of prefetched pages, until it stops prefetching completely. Once these demand and

prefetched pages arrive, they are placed into the swap cache. Their swap entries in remote memory are then freed.

The kernel uses an LRU algorithm to evict pages. Evicting a page *unmaps* it and pushes it into the swap cache. When memory runs low, the kernel releases existing pages from the swap cache to make room for newly fetched pages. Clean pages can be removed right away and dirty pages need to be written back. To write back a page, the swap system must first allocate a swap entry using a free-list-based allocation algorithm for the page. Finally, an I/O request is generated and the page is written into the entry via RDMA.

In each remote access, extensive locking is needed for swap entry allocation—shared allocation metadata (*e.g.*, free list) must be protected (similarly to how a memory allocator protects its metadata) when multiple applications/threads request swap entries simultaneously. Although there are active efforts [46, 44] in the kernel development community to optimize swap entry allocation, their performance and scalability is unsatisfactory for cloud workloads (see §5.1).

3 Motivating Performance Study

To understand the impact of interference, we conducted a study with a set of widely-deployed applications including Apache Spark [104], XGBoost [21] (*i.e.*, a popular ML library), Snappy [36] (*i.e.*, Google’s fast compressor/decompressor), as well as Memcached [4]. Spark is a managed application running on the JVM, while the other three are native applications. They cover a spectrum of cloud workloads from data storage through analytics to ML.

We ran these programs, individually vs. together, on a machine with two Xeon(R) Gold 6252 processors, running Linux 5.5. Another machine with two Xeon(R) CPU E5-2640 v3 processors and 128GB memory was used for remote memory. Each machine was equipped with a 40 Gbps Mellanox ConnectX-3 InfiniBand adapter and interconnected by one Mellanox 100 Gbps InfiniBand switch. Using `cgroup`, the same amounts of CPU and local memory resources were given to each application throughout the experiments. RDMA bandwidth was *not* saturated for both application individual runs and co-runs. The amount of local memory configured for each application was 25% of its working set.

Performance Interference and Degradation. To understand the overall performance degradation and how it changes with different applications, we used two Spark applications: PageRank (SPR) and Logistic Regression (SLR). Figure 2 reports each application’s performance degradation when co-running with other applications compared to running alone. The blue/orange bars show the slowdowns when the three native applications co-run with SPR/SLR. Clearly, co-running applications significantly reduces each application’s performance. We observed an overall $2.24/3.75\times$ slowdown when native applications co-run with SPR/SLR. SLR persists a large RDD in memory and keeps swapping in/out different

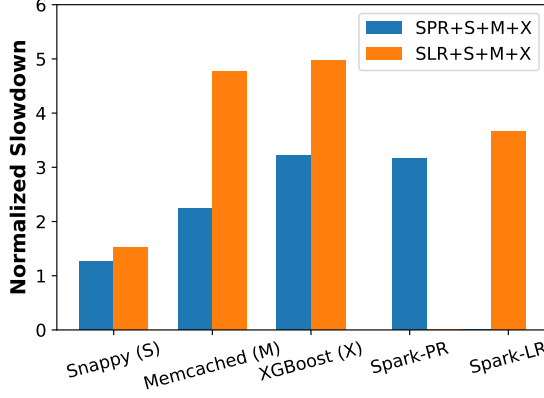


Figure 2: Slowdowns of co-running applications compared to running each individually.

parts of the RDD, while SPR holds much of its graph data in local memory and thus does not swap as much as SLR.

Another observation is that the impact of interference differs significantly for different applications. Applications that generate high swap throughputs aggressively invade swap and RDMA resources of other applications. In our experiments, Memcached, XGBoost, and Spark all need frequent swaps. However, Spark runs many more threads (>90 application and runtime threads) than Memcached (4) and XGBoost (16), resulting in a much higher swap throughput. As such, Spark takes disproportionately more resources, leading to severe degradation for Memcached and XGBoost.

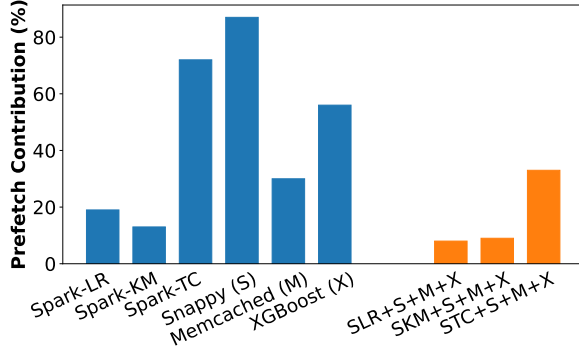


Figure 3: Prefetching contribution of LEAP: the percentage of page faults served by LEAP-prefetched pages (%).

Reduced Prefetching Effectiveness. Sharing the same prefetching policy reduces the prefetching effectiveness when multiple applications co-run. Figure 3 reports *prefetching contribution*—the percentage of page faults served by prefetched pages—the higher the better; if a prefetched page is never used, prefetching it would only incur overhead. We used LEAP [70] as our prefetcher. The left six bars report such percentages for the applications running individually. When applications co-run, the rightmost three bars report

the average percentages across applications. As shown, co-running dramatically reduces the contribution.

Note that LEAP [70] uses a majority-vote algorithm to identify patterns across multiple applications. However, when applications that exhibit drastically different behaviors co-run, LEAP cannot adapt its prefetching mechanism and policy to each application. Furthermore, LEAP is an aggressive prefetcher—even if LEAP does not find any pattern, it always prefetches a number of contiguous pages. However, aggressive prefetching for applications such as Spark with garbage collection (GC) is ineffective—e.g., prefetching for a GC thread has zero benefit and only incurs overhead. Detailed evaluation of prefetching can be found in §6.3.

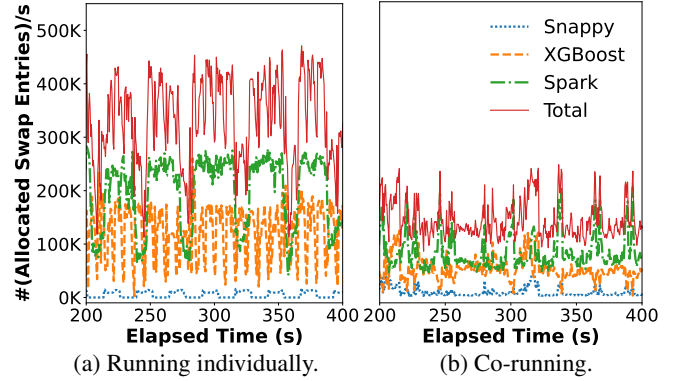


Figure 4: Swap entry allocation throughput when applications run individually (a) and together (b).

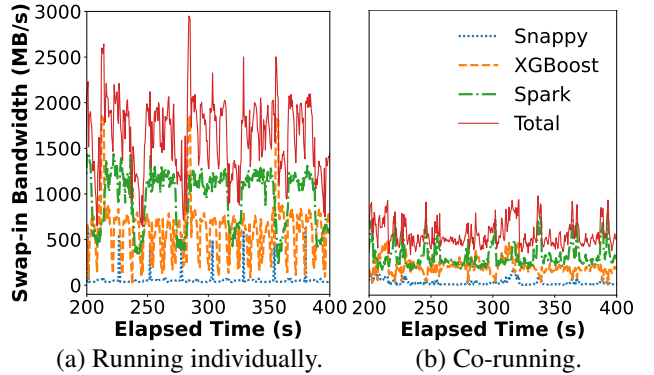


Figure 5: RDMA swap-in bandwidth when applications run individually (a) and together (b).

Lock Contention. We observed severe lock contention in the swap system when applications co-run, particularly at swap entry allocation associated with each swap-out.

We experimented with Spark (Logistic Regression), XGBoost, and Snappy. Our results show that in windows of frequent remote accesses, co-running applications can spend up to **70%** of the window time on obtaining swap entries. Lock contention leads to significantly reduced swap-entry al-

location throughput, reported in Figure 4. The total lines in Figure 4(a) and (b) show the total throughput (*i.e.*, the sum of each application’s allocation throughput). The co-running throughput (b) is drastically reduced compared to the individual run’s throughput (a) (*i.e.*, $\sim 450\text{Kps}$ to $\sim 200\text{Kps}$).

Reduced RDMA Utilization. Figure 5 compares the RDMA read bandwidth (for swap-ins) when applications run individually and together. Similarly, the total line represents the sum of each application’s RDMA bandwidth. The total RDMA utilization is constantly below $\sim 1000\text{MBps}$ in Figure 5(b), which is $3.28\times$ lower than that in Figure 5(a) due to various issues (*e.g.*, locking, reduced prefetching, *etc.*). The RDMA write bandwidth degrades by an overall of $2.80\times$.

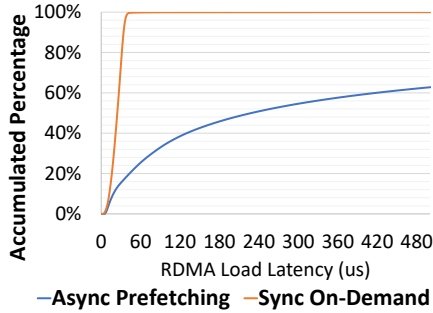


Figure 6: Latency of prefetching and on-demand swapping.

Demand v.s. Prefetching Interference. Optimizations such as Fastswap [8] improve swap performance by dividing the RDMA queue pairs (QP) into sync and async. The high-priority synchronous QP is used for demand swaps, while the low priority async QP is used for prefetching requests. This separation reduces head-of-line blocking incurred by prefetching. However, when applications co-run, this design adds a delay for prefetching. Figure 6 depicts the CDF of the latency of RDMA packets from demand and prefetching requests, when the four applications co-run on Fastswap (with the LEAP prefetching algorithm). As shown, 99% of the on-demand requests are served within $40\mu\text{s}$. However, the latency of 36.9% of prefetching requests is longer than $512\mu\text{s}$ and it can reach up to 52ks ! This creates two problems. First, long latency renders prefetched pages useless because prefetching is meant to load pages to be used soon. Our profiling shows that among the prefetched pages that are actually accessed by the application, 90% are accessed within $70\mu\text{s}$, indicating that $\sim 70\%$ of the pages prefetched return too late. Second, prefetching such pages wastes swap cache space and RDMA bandwidth. These problems motivate our two-dimensional RDMA scheduling (§5.3).

Takeaway. The root cause of performance degradation is that multiple applications, whose resource needs and swap behaviors are widely apart, all run on a global swap system with the same allocator and prefetcher. Table 1 summarizes these problems, their performance impact, and our solutions.

4 Swap System Isolation

Canvas extends `cgroup` for users to specify size constraints for swap partition (*i.e.*, remote memory), swap cache, and RDMA bandwidth. We discuss the kernel support to enforce these new constraints, laying a foundation for adaptive optimizations discussed in §5.

Swap Partition Isolation. Remote memory is managed via a swap partition interface, which consists of a set of swap entries. In Linux, each entry has a 4KB size and a swap partition is shared by all applications. If there are multiple available swap partitions, they are used in a *sequential manner* according to their priorities. As a result, data of different applications are mixed and stored in arbitrary locations.

Canvas separates remote memory of each `cgroup` to isolate capacity and performance. The user creates a `cgroup` to set a size limit of remote memory for an application. Canvas allocates remote memory in a demand-driven manner—upon a pressure in local memory, Canvas allocates remote memory and registers it as a RDMA buffer. Canvas enables per-`cgroup` swap partitions by creating a swap partition interface and attaching it to each `cgroup`. For each `cgroup`, a separate swap-entry manager is used for allocating and freeing swap entries. Swap entry allocation can now be charged to the `cgroup`, which controls how much remote memory each application can use. Our adaptive swap entry allocation algorithm is discussed in §5.1.

Canvas explicitly enables a private swap cache for each `cgroup` (a default value of 32MB), whose size is charged to the *memory budget* specified in the `cgroup`. As a result, the size of an application’s swap cache changes in response to its own memory usage, without affecting other applications.

For each demand swap-in, Canvas first checks the `mapcount` of the page, which indicates how many processes this page has been mapped to before. If the page belongs only to one process, it is placed in its private swap cache. Otherwise, it has to be placed in a global swap cache (discussed shortly). To release pages (*e.g.*, when the application’s working set increases, pushing the boundary of the swap cache), Canvas scans the swap cache’s page list, releasing a batch of pages to shrink the cache.

RDMA Bandwidth Isolation. For each `cgroup`, Canvas isolates RDMA bandwidth with a set of *virtual* RDMA queue pairs (VQPs) and a centralized packet scheduler. Users can set the swap-in/swap-out RDMA bandwidth of a `cgroup` with our extended interface. Our RDMA scheduler works in two dimensions. The *first dimension* schedules packets across applications, while the *second dimension* schedules on a per-application basis—each `cgroup` has its *sub-scheduler* that schedules packets that belong to the `cgroup` between demand swapping and prefetching.

VQPs are high-level interfaces, implemented with lock-free linked lists. Each `cgroup` pushes its requests to the head of its VQP, while the scheduler pops requests from their tails. At the low level, our scheduler maintains three

Problem Description	Performance Impact	Canvas’s Solution
Unlimited use of swap and RDMA resources	Apps generating higher swap thrupt use disproportionately more resources	Holistic isolation of swap system RDMA isolation and scheduling (§4, §5.3)
Lock conten. at swap entry alloc.	Reduced swap-out throughput	(1) Swap parti. isolation (§4); (2) adaptive entry alloc. (§5.1)
Single low-level prefetcher	Increased fault-handling latency	Two-tier adaptive prefetching (§5.2)
prefetching v.s. demand interfere	Increased fault-handling latency	Two-dimensional RDMA scheduling (§5.3)

Table 1: Summary of major issues and Canvas’s solution.

physical queue pairs (PQP) per core, for *demand swap-in*, *prefetching*, and *swap-out*, respectively. The scheduler polls all VQPs and forwards packets to the corresponding PQPs, using a *two-dimensional* scheduling algorithm (see §5.3).

Handling of Shared Pages. Processes can share pages due to shared libraries or memory regions. These pages cannot go to any private swap cache. Canvas maintains a global swap partition and cache for shared pages. When a page is evicted and unmapped, Canvas checks its `mapcount` and adds it to the global swap cache if the page is shared between different processes. All pages in the global swap cache will be eventually swapped out to the global partition using the original lock-based allocation algorithm. Conversely, pages swapped in (and prefetched) from the global swap partition are all placed into the global swap cache. Given that the number of shared pages is much smaller than process-private pages, using locks in a normal way would not incur a large overhead. We cannot charge applications’ `cgroups` for pages in the global swap cache, because which process(es) share these pages is unknown before they get mapped into processes’ address spaces. Canvas allows users to create a special `cgroup`, named `cgroup-shared`, to limit the size of the global swap cache/partition.

5 Isolation-Enabled Swap Optimizations

On top of the isolated swap system, we develop three optimizations, which dynamically adapt their strategies to each application’s resource patterns and semantics.

5.1 Adaptive Swap Entry Allocation

As discussed in §3, swap entry allocation suffers from severe lock contention under frequent remote accesses—allocation is needed at every swap-out. Creating a per-application swap partition mitigates the problem to a certain degree. However, applications like Spark often create a great number of processing threads, still incurring significant locking overhead.

To further reduce contention, we develop a novel swap entry allocator that adapts allocation strategies in response to each application’s own memory access/usage. Our first idea is to enable a *one-to-one* mapping between pages and swap entries. At the first time a page is swapped out, we allocate a new swap entry using the original (lock-protected) algorithm. Once the entry is allocated, Canvas writes the entry ID into the page data (as a piece of 8-byte metadata). This ID remains on the page throughout its life span. As a result,

subsequent swap-outs of the page can write data directly into the entry corresponding to this ID. We pay the locking overhead *only once* for each page at its first swap-out.

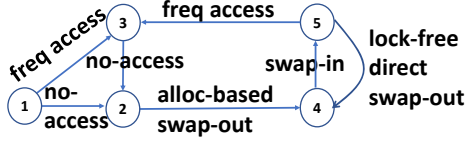
This approach requires a swap entry to be reserved for each page. An entry reserved for one page cannot be used to store any other pages due to the permanent one-to-one mapping. This may unnecessarily increase the remote memory usage. For example, modern cloud applications exhibit strong epochal behaviors. Under the original allocator, swap entries for pages accessed in one epoch can be reused for those in another epoch. Under this approach, however, all pages in all epochs must have their dedicated swap entries throughout the execution, which can lead to an order-of-magnitude increase in remote memory usage.

Our key insight is: we should trade off *space for time* if an application has much available swap space, but *time for space* when its space limit is about to be reached. As such, when the remote memory usage is about to reach the limit specified in `cgroup` (*i.e.*, 75% in our experiments), Canvas starts removing reservations to save space. The next question is which pages we should consider first as our candidates for reservation removal. Our idea is that we should first consider “hot pages” that always stay in local memory and are rarely swapped. Our observation is: hot pages (*i.e.*, data on such pages are frequently accessed) are likely to stay in local memory for a long time; hence, locking overhead is less relevant for them. On the contrary, “cold” pages whose accesses are *spotty* are more likely to be swapped in/out and hence swap efficiency is critical. Here “hot” and “cold” pages are relatively defined as they are specific to execution stages—a cold page swapped out in a previous stage can be swapped in and become hot in a new stage.

To this end, we develop an *adaptive allocator*. Canvas starts an execution by reserving swap entries for *all* pages to minimize lock contention. Reservation removal begins when remote-memory pressure is detected. Canvas adaptively removes reservations for hot pages. We detect hot pages *for each application* by periodically scanning the application’s *LRU active list*—pages recently accessed are close to the head of the active list. Each scan identifies a set of pages from the head of the list; a page is considered “hot” if it appears in a consecutive number of sets recently identified.

Removing reservation for a hot page can be done by (1) removing the entry ID from the page metadata and (2) freeing

its reserved swap entry in remote memory, adding the entry back to the free list. Once a hot page becomes cold and gets evicted, it does not have a reservation any more, and hence, it goes through the original (lock-protected) allocation algorithm to obtain an entry. In this case, the page receives a new swap entry and remembers this new ID in its metadata.



1. Init
2. Cold page w/o swap entry ID
3. Hot page, swap-entry ID removed
4. In remote memory
5. Cold page with entry ID

Figure 7: FSM describing our page management when remote-memory pressure is detected.

Figure 7 shows the page state machine, which describes the page handling logic. A cold page (to be evicted) can be in one of the two states: state 2 and state 5. A page comes to state 2 if it is (1) a brand new page that has never been swapped out or (2) previously a hot page but has not been accessed for long. Once it reaches state 2, the page does not have a reserved swap entry ID and hence, swapping out this page goes through the normal allocation path. In the case of swap-in (state 5), the swap entry ID is already remembered on the page. The next swap-out will directly use this entry and be lock-free. If the page becomes hot (from state 5 to 3), Canvas removes the entry ID and releases the entry reservation. The entry is then added back to the free list.

Recent Kernel Development. As an optimization in Linux 5.5, the kernel keeps swap entries for clean pages—when clean pages are evicted, they do not need to be written back if their swap entries are not released for other allocations. Once a page becomes dirty, its swap entry must be immediately released. Clearly, this approach works for read-intensive applications where most pages are clean, but not for write-intensive workloads such as Spark. We tried various entry-keeping thresholds (*i.e.*, entry keeping starts when the percentage of available swap entries exceeds this threshold) between 25% and 75%, and saw only marginal performance differences (<5%) across our programs.

We have closely followed the kernel development since the release of Linux 5.5 and found two recent patches related to our approach. These two patches, submitted by Intel and merged into the kernel at 5.8, also attempt to optimize locking overhead at swap entry allocation. The idea of the first patch [46] is using fine-grained locking—dividing swap entries into *clusters* and assigning each core a random cluster upon an allocation request. The second patch [44] performs batch entry allocation by scanning more swap entries while holding the lock to make each batch larger. Note that

our adaptive allocation algorithm solves a much bigger problem than these patches—Canvas *avoids* allocating entries for most swap-outs, while these patches reduce the overhead of locking for each allocation. As such, Canvas is completely lock-free for reserved entries while these patches must still go through the allocation path, requiring locking if multiple cores are assigned the same cluster (*i.e.*, core collision).

In fact, the probability of collision increases quickly with the number of cores. As reported in §5.1, the allocation performance of these patches degrades super-linearly when the number of cores exceeds 24. Another major drawback is that none of these patches build on isolated swap partitions. Lack of swap partition isolation makes applications search for swap entries globally, which can still result in interference—applications such as Spark can quickly saturate these clusters with all its executor threads, making other applications wait before they can obtain the locks. By reserving entries, our algorithm significantly reduces the number of entry allocation requests (due to entry reusing) and the cost of each allocation (due to reduced lock contention).

As the kernel is fast evolving and our latest InfiniBand driver is only compatible with Linux 5.5, we compared the allocation performance between the latest Linux 5.14 and Canvas over RAMDisk. As detailed in §6.3.1, the Linux 5.14 has severe scalability problems. With 48 cores, our algorithm outperforms Linux 5.14’s entry allocator (that uses [46, 44]) by 13×.

5.2 Two-tier Adaptive Prefetching

Problems with Current Prefetchers. Current prefetchers all focus on low-level (streaming or strided) access patterns. While such patterns exist widely in native array-based programs, applications written in high-level languages such as Python and Java are dominated by reference-based data structures—operations over such data structures involve large amounts of pointer chasing, making it hard for current prefetchers to identify clear patterns.

Furthermore, cloud applications such as Spark are heavily multi-threaded. Modern language runtimes such as the JVM runs an additional set of auxiliary threads, *e.g.*, for GC or JIT compilation. How these user-level threads map to kernel threads is often implemented differently in different runtimes. Consequently, kernel prefetchers such as LEAP [70] cannot distinguish patterns from different threads.

To develop an adaptive prefetcher, Canvas employs a two-tier design, illustrated in Figure 8. At the low (kernel) tier, Canvas uses an existing kernel prefetcher that prefetches data for each application into its own private swap cache (unless data comes from the global swap partition). A kernel prefetcher is extremely efficient and can already cover a range of (array-based) applications. For applications whose accesses are too complex for the kernel prefetcher to handle, we forward the addresses up to the application level, letting

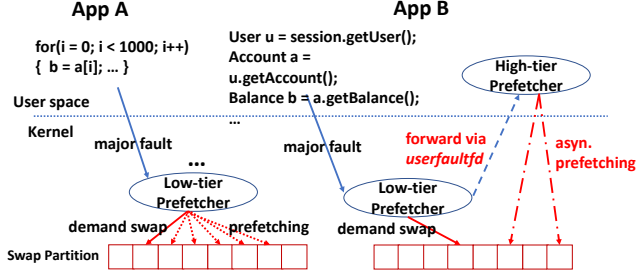


Figure 8: Canvas’s two-tier prefetcher: App A is an array-based program while B is a modern web application that uses reference-based data structures. The low-tier prefetcher successfully prefetches pages for A, but not for B. Hence, Canvas forwards the addresses up to B’s high-tier prefetcher.

the application/runtime analyze semantic access patterns at the level of threads, references, arrays, *etc.*

Prefetching Logic. In Canvas, we adopt the sync/async separation design in Fastswap [8], which prevents head-of-line blocking. As stated earlier, we use three PQPs per core, one for swap-out, one for (sync) demand swap-in, and one for (async) prefetching. Canvas polls for completions of critical (demand) operations, while configuring *interrupt completions* for asynchronous prefetches.

Canvas determines whether to use an application-tier prefetcher based on *how successful kernel-tier prefetching is*. If the number of pages prefetched for an application is lower than a threshold at the most recent N ($=3$ in our evaluation) faults consecutively, Canvas starts forwarding the faulty addresses up to the application-tier prefetcher (discussed shortly) although the kernel-tier prefetcher is still used as the first-line prefetcher.

Canvas stops forwarding whenever the kernel-tier prefetcher becomes effective again. Our key insight is: the kernel-tier prefetcher is efficient without needing additional compute resources (as it uses the same core as the faulting thread), while the application-tier prefetcher needs extra compute resources to run. As such, we disable application-tier prefetchers as long as the kernel-tier prefetcher is effective. To pass a faulty address to the application, we modify the kernel’s `userfaultfd` interface, allowing applications to handle faults at the user space. Our modification makes the kernel forward the faulty address only if the kernel’s prefetcher continuously fails to prefetch pages.

Runtime Support for Application-tier Prefetching. A major challenge is how to develop application-tier prefetchers. On the one hand, application-tier prefetchers should conduct prefetching based on *application semantics*, of which the kernel is unaware. On the other hand, prefetching is such a low-level activity that application developers may not be familiar with; understanding memory access patterns and developing prefetchers can be a daunting task for them.

Our insight is: applications that benefit from application-tier prefetching are mostly written in high-level languages and run on a managed runtime such as the JVM. Inspired by previous work on using language runtime to solve memory efficiency problems for data analytics applications [77, 75, 78, 76, 69], Canvas currently supports application-tier prefetching for the JVM as a platform. However its support could be easily extended to other managed runtimes for high-level languages like Go and C#. Leveraging language runtime solves both problems discussed above—it has access to semantic information such as how objects are connected and the number of application threads; furthermore, the burden of developing an application-tier prefetcher is shifted from application developers to runtime developers. Thus, it is not necessary to supply a custom application-tier prefetcher per application, but define it once for each language runtime.

In this work, we develop an application-tier prefetcher in Oracle’s OpenJDK as a proof-of-concept. It works for all (Java, Scala, Python, *etc.*) programs that run on the JVM. Our JVM-based prefetcher considers two *semantic patterns*: (1) *reference-based* (*i.e.*, accessing an object brings in pages containing objects referenced by this object) and (2) *thread-based* (*i.e.*, accesses from different application threads are separately analyzed to find patterns).

For (1), we modify the JVM to add support that can quickly find, from a faulty address, *the object* in which the address falls. Next, Canvas uses a garbage collection like approach that follows a chain of pointers (on a coarse-grained reference summary graph constructed by GC and write barrier) to identify pages referenced by the faulty object. Our implementation chases pointers for 3 hops. Canvas then notifies the kernel, via a system call `asyn_prefetch`, of the prefetching requests for these pages. Finally, RDMA fetches them into the application’s swap cache.

For (2), we leverage the JVM’s user-kernel thread map. For each faulty address, Canvas additionally forwards the thread information (*i.e.*, pid) to the JVM, which consults the map to filter out non-application (*e.g.*, GC, compilation, *etc.*) threads and segregate addresses based on Java threads (as opposed to kernel threads). Segregated addresses allow us to analyze (sequential/strided) patterns on a per-thread basis (using LEAP’s majority-vote algorithm [70]). Once patterns are found, the prefetcher sends the prefetching requests to the kernel via `asyn_prefetch`.

Policy. To improve effectiveness, the JVM uses a search tree to record information about large arrays. Upon the allocation of an array whose size exceeds a threshold (*i.e.*, 1MB in our experiments), the JVM remembers its starting address and size into the tree. The JVM runs a daemon prefetching thread. Once receiving a sequence of faulty addresses, we determine which semantic pattern to use based on *how many application threads are running* and *whether the faulty addresses fall into a large array*. If there are many threads and the faulty addresses fall into arrays, the JVM uses (2) to

find per-thread patterns. If either condition does not hold, the JVM uses (1) to prefetch based on references. For native applications, we only enable (2), as we observed that our native programs do not use many deep data structures.

5.3 Two-Dimensional RDMA Scheduling

To provide predictable performance for applications sharing RDMA resources, our RDMA scheduling algorithm should provide four properties: (1) weighted fair bandwidth sharing [16, 28] across applications; (2) high overall utilization; (3) treating demand and prefetching requests with different priorities; and (4) timely handling of prefetching requests.

Canvas performs two-dimensional scheduling by extending existing techniques. Canvas uses max-min fair scheduling to assign bandwidth across applications, and priority-based scheduling with *timeliness* to schedule prefetching and demand requests within each application. Although these scheduling techniques are not new themselves, Canvas combines them in a unique way to solve the interference problem. Canvas maintains three PQPs on each core, respectively, for swap-outs, demand swap-ins, and prefetching swap-ins. Swap-outs are only subject to fair scheduling while swap-ins are to both fair and priority-based scheduling.

Vertical: Fair Scheduling. Under max-min fairness, each application receives a fair share of bandwidth. If there is extra bandwidth, we give it to the applications in the reverse order of their bandwidth demand until bandwidth is saturated. The high overall utilization of bandwidth is achieved by redistributing unconsumed bandwidth proportionally to the weights of unsatisfied applications. Canvas implements weighted fair queuing with virtual clock [80, 28, 105].

Horizontal: Priority Scheduling with Timeliness. Within each `cgroup`, Canvas schedules demand requests with a higher priority than prefetching requests. However, this could lead to long latency for prefetching requests. To bound the latency of prefetching, our scheduler employs a history-based heuristic algorithm to identify and drop outdated prefetching requests. In particular, Canvas maintains the *timeliness distribution* of prefetched pages per `cgroup`. Timeliness is a metric that measures the time between a page being prefetched and accessed. We attach a timestamp to each request when pushing it into a VQP. The scheduler maintains packets statistics on-the-fly to estimate the round-trip latency and arrival time of each prefetching request. Requests are dropped if the estimated arrival time exceeds the estimated timeliness threshold. Special care must be taken to drop prefetching requests. Before issuing a prefetching request, the kernel creates a page in the swap cache and sets up its corresponding PTE. The page is left in a *locked* state until its data comes back. However, a thread that accesses an address falling into the page may find this locked page in the swap cache and block on it. Dropping prefetching requests may cause the thread to hang. To solve the problem, we de-

test threads that block on prefetching requests for too long and generate new *demand requests* for them.

6 Evaluation

It took us 17 months to implement Canvas in Linux 5.5. The application-tier prefetcher was implemented in OpenJDK 12.

Spark Applications	Dataset	Size
Spark PageRank (SPR)	Wikipedia Polish [5]	1GB
Spark KMeans (SKM)	Wikipedia English [5]	3GB
Spark Logistic regression (LR)	Wikipedia English [5]	3GB
Spark SkewedGroupby (SSG)	Wikipedia Polish [5]	1GB
MLlib Bayes Classifiers (MLB)	KDD [3]	10GB
GraphX ConnectedComponents (GCC)	Wikipedia English [5]	3GB
GraphX PageRank (GPG)	Wikipedia English [5]	3GB
GraphX Single Source Shortest Path (GSSSP)	Spark synthetic	2M vertices
GraphX Triangle Counting (GTC)	A synthetic dataset in Spark	1.5M edges, 384K vertices
Native Applications		
XGBoost	HIGGS[11]	22M instances
Snappy	enwik9 [1]	16GB
Memcached	YCSB[24]	10M records, 45M gets, 5M sets

Table 2: Programs and their datasets.

Setup. We included a variety of real-world applications in our experiments, including nine Spark applications as well as three native applications: XGBoost [21], Snappy [36], and Memcached [4]. Spark, Memcached, and XGBoost are multi-threaded while Snappy is single-threaded. The Spark applications span several popular libraries such as GraphX and MLlib. The performance of their individual runs on remote memory is consistent with what was reported in prior work [100, 8], and omitted due to space constraints.

We co-ran different combinations of programs. The same application in different combinations receives the same amount of local (CPU and memory) resources. To simplify performance analysis, we let each combination of applications we co-run contain one Spark application and the three native programs (which consume less resources than Spark). These experiments were conducted on two machines, one used to execute applications and a second to provide remote memory. The configurations of these machines was reported earlier in §3. We carefully configured Linux with the following configuration to achieve the best performance for Linux: (1) SSD-like swap model, (2) per-VMA prefetching policy, and (3) cluster-based swap entry allocation. We disabled hyper-threads, CPU C-states, dynamic CPU frequency scaling, and transparent huge pages.

We limited the amounts of CPU resources for Spark, XGBoost, Memcached, and Snappy to be 24, 16, 4, and 1 core(s). For local memory, we used two ratios: 50% and 25%, meaning each application has 50/25% of its working set locally. When using Canvas, we limited the sizes of swap partitions in such a way that for each application the total size of its swap partition and assigned local memory is slightly larger than its working set. In doing so, each application has just enough (local and remote) memory to run and reservation cancellation (§5.1) is triggered in all executions.

The swap cache size for each application starts at 32MB and changes dynamically. The global swap cache size (con-

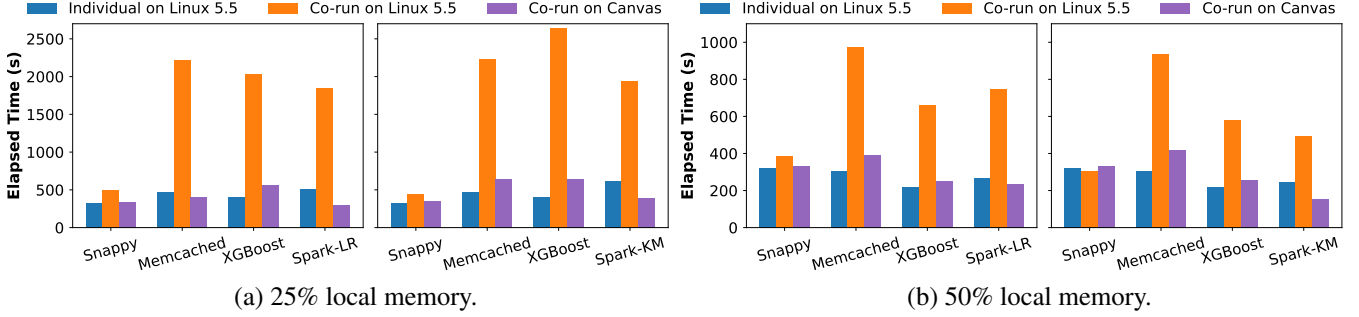


Figure 9: Overall performance under 25% and 50% local memory when native programs co-run Spark LR and Spark KM.

figured by `cgroup-share`) was also set to 32MB. The RDMA bandwidth assignments were proportional to their swap partition assignments.

6.1 Overall Performance

We first demonstrate the overall performance when applications co-run together under Canvas. Each experiment ran the same set of three native programs with a different Spark application. Due to space constraints, we only report the performance comparisons for the runs that contain Spark LR (Figure 9(a)) and Spark KM (Figure 9(b)).

The three bars in each group represent an application’s performance when running alone on Linux 5.5, co-running with other applications on Linux 5.5, and co-running on Canvas (with all optimizations enabled). Across all experiments with the nine Spark applications, Canvas improves applications’ co-run performance by up to $6.1\times$ (average $2.35\times$) and up to $3.6\times$ (average $1.33\times$) under the two memory configurations. Canvas enables Spark to even outperform their individual runs due to the optimizations that could also improve single-application performance.

6.2 Isolation Reduces Degradation and Variation

Here we focus on understanding the effectiveness of isolation alone. We used a variant of Canvas with the isolated swap system and RDMA bandwidth (*i.e.*, vertical scheduling between applications) but without the adaptive swap-entry scheduler, application-tier prefetcher, and horizontal RDMA scheduling (demand swapping vs. prefetching).

Degradation Reduction. We ran the three native applications with each of Spark-LR, Spark-KM, and Spark-PR. As shown in Figure 10, isolation reduces the running time by 1.05 - $5.65\times$, with an average of $2.22\times$ for native applications. Isolation is particularly useful for applications that do not have many threads but need to frequently access remote memory, such as Memcached, which has 4 threads and cannot compete for resources with Spark with more than 90 (application and runtime) threads. As such, its performance is improved by $2.48\times$ with dedicated swap resources. Isolation improves the average RDMA utilization by $2.8\times$

from **670.9MB/s** to **1878.4MB/s**, making the peak bandwidth reach 4493.5MB/s (*i.e.*, RDMA bandwidth saturated).

Slowdown	Mean	Min	Max	σ
Snappy	1.03 / 1.23	1.03 / 1.09	1.05 / 1.53	0.01 / 0.15
Memcached	1.09 / 2.74	1.03 / 1.35	1.17 / 5.27	0.06 / 1.58
XGBoost	1.10 / 3.31	1.06 / 1.90	1.17 / 6.44	0.04 / 1.57
Overall	1.08 / 2.42	1.03 / 1.09	1.17 / 6.44	0.05 / 1.53

Table 3: Performance statistics of three native applications when co-running with each of the nine Spark applications under 25% local memory setting (Canvas / Linux 5.5).

Variation Reduction. One significant impact of interference is performance variation—the same application has drastically different performance when co-running with different applications (as shown in Figure 2). To demonstrate our benefits, we co-ran the three native applications with each of the nine Spark applications listed in Table 2, which cover a wide spectrum of computation and memory access behaviors. Table 3 reports various statistics of their performance including the mean, minimum, maximum, and standard deviation of their execution times. Clearly, the performance of the three programs is much more stable under Canvas than Linux—variations are reduced by an overall of $31\times$.

6.3 Effectiveness of Adaptive Optimizations

Our goal here is to understand the benefit of each swap optimization *on top of the isolated swap system*. We evaluated each optimization’s contribution by turning it on/off while leaving all other optimizations enabled.

6.3.1 Adaptive Swap Entry Allocator

Isolation already reduces lock contention at swap entry allocation because each process has its own swap entry manager. However, for multi-threaded applications such as Spark, their processing threads still have to go through the locking process. In this subsection, we focus on Spark due to its extensive use of threads. Figure 11 shows the performance of Spark LR and Spark KM when they each co-run with the other three native programs. On average, our adaptive allocation enables an *additional* boost of $1.53\times$ for Spark LR

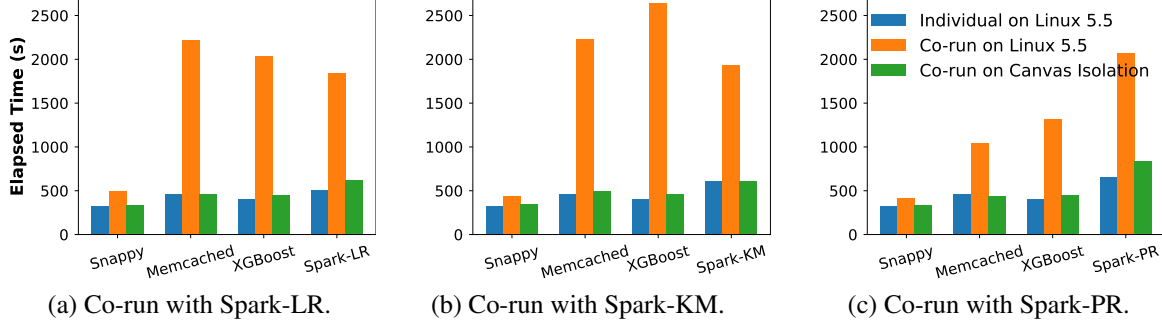


Figure 10: Native applications co-run with different Spark applications under isolation *alone*.

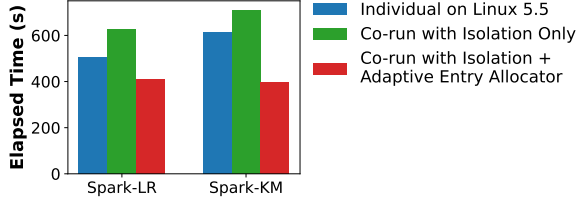


Figure 11: Benefit of adaptive swap entry allocation for two Spark applications. Compared are the times of the application running individually on Linux 5.5, co-running on Canvas with adaptive entry allocation disabled, and enabled.

and $1.56\times$ for Spark KM, making both of them outperform their individual runs.

Thruput. (KPages/s)	Linux 5.5	Canvas w/o adap. alloc.	Canvas w/
Avg. Spark apps	98	164	295
Avg. all apps	185	309	468

Table 4: Swap-out throughput with and without adaptive swap-entry allocation.

Table 4 reports the swap-out throughput when four applications co-run. As shown, isolation improves the throughput by $1.67\times$ while adaptive allocation provides an additional improvement of $1.51\times$. This benefit is obtained after applying all optimizations in Linux 5.5.

Comparison with Linux 5.14 on RAMDisk. We compared our adaptive allocation algorithm with the allocator in Linux 5.14 that uses the patches [46, 44] by running Memcached with varying (8 – 48) cores. Since our RDMA is incompatible with any post-5.5 kernel version, we could not run Linux 5.14 directly on remote memory. Instead, we ran this experiment on RAMDisk with a focus on allocation performance.

As Figure 12(a) shows, our adaptive entry reservation algorithm reduces the allocation rate by an order of magnitude compared to Linux 5.14. Figure 12(b) compares our algorithm with Linux 5.14 on per-entry allocation time. As shown, the optimization in [46, 44] is unscalable—as the number of cores increases, the per-entry allocation cost in-

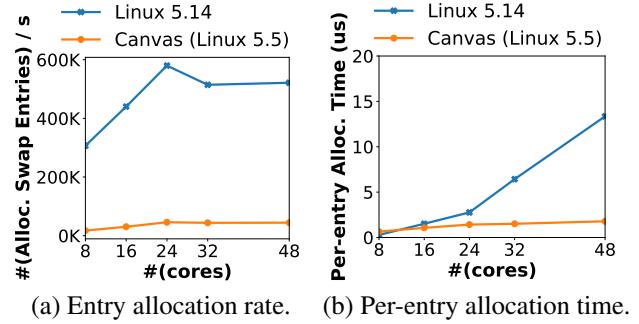


Figure 12: Swap entry allocation comparison between Canvas and the allocation algorithm in Linux 5.14 on RAMDisk.

creases significantly. In fact, the allocation cost grows super-linearly after 24 cores due to core collision. On the contrary, Canvas’s per-entry allocation cost remains low and stable.

6.3.2 Prefetching Effectiveness

To demonstrate the benefit of application-tier prefetching, our baseline here is the kernel’s default prefetcher on the isolated swap system with the other two optimizations (*i.e.*, adaptive swap allocator and RMDA scheduling) *enabled*. Since application-tier prefetching is designed primarily for high-level languages, here we focus on Spark applications.

Time. We compare the running time for three Spark applications: LR, KM, and TC, between the kernel’s prefetcher over Canvas’s isolated swap system and Canvas’s two-tier prefetcher, when each Spark application co-runs with the three native applications under the 25% local memory configuration. Application-tier prefetching brings **33%**, **17%**, and **19%** additional performance benefits on top of the kernel prefetching with the isolated swap system. All the three Spark applications benefit from the thread-level pattern analysis while LR and KM have seen 5-9% contributions from using the reference-based pattern. The thread-level pattern analysis we added for native programs brings a 5% and 11% improvement for Memcached and XGBoost.

We have also run LEAP [70], a prefetcher that aggressively prefetches a number of contiguous pages if it cannot

find any pattern. This approach may work for native applications because these applications access arrays; hence, the contiguous pages aggressively prefetched are likely to be useful for array accesses. However, it works poorly for high-level language applications such as Spark, which uses deep data structures and runs graph-traversal GC tasks (for which neither sequential nor strided prefetching is useful). Aggressively prefetching useless pages wastes the RDMA bandwidth and the swap cache space. LEAP slows down Spark by $1.4\times$, compared to the kernel’s default prefetcher.

Contribution	Spark-LR	Spark-KM	Spark-TC
LEAP	23.4%	25.8%	42.2%
Kernel	63.3%	68.0%	65.9%
Canvas Two-tier	79.2%	79.3%	75.3%
Accuracy	Spark-LR	Spark-KM	Spark-TC
LEAP	16.8%	17.2%	35.9%
Kernel	95.6%	96.4%	93.9%
Canvas Two-tier	94.3%	94.8%	94.9%

Table 5: Prefetching contribution and accuracy when different Spark workloads co-run with native applications.

Prefetching Contribution and Accuracy. Table 5 compares prefetching *contribution* and *accuracy* for the three Spark applications when each of them co-runs with the same three native applications. Contribution is defined as a ratio between the number of page faults hitting on the swap cache and the total number of page faults (including both cache hits and demand swap-ins). Accuracy is defined as a ratio between the number of page faults hitting on the swap cache and the total number of prefetches. Clearly, contribution has a strong correlation with performance while accuracy measures the pattern recognition ability of a prefetcher. For example, for a conservative prefetcher that prefetches pages only if a pattern can be clearly identified, it can have a high accuracy (*i.e.*, prefetched pages are all useful) but a low contribution (*i.e.*, the number of prefetches is small).

Here we report prefetching contribution and accuracy for three prefetchers: LEAP (on our isolated swap system), the kernel prefetcher (also on our isolated swap system), and Canvas’s two-tier prefetcher. Among the three prefetchers, LEAP has the lowest accuracy and contribution because it is an aggressive prefetcher. First, LEAP keeps prefetching pages even when it cannot detect any patterns, which greatly reduces the prefetching accuracy. Second, due to the limited swap cache, the useless pages prefetched can cause previously prefetched pages to be released before they are accessed. As a result, the aggressiveness also harms the contribution. The kernel prefetcher and Canvas have comparable accuracy because the kernel prefetcher is much more conservative than LEAP. It stops prefetching when no clear pattern can be observed. However, it has lower contribution than our two-tier prefetcher since Canvas can prefetch more useful pages using semantic information.

6.3.3 RDMA Scheduling

We evaluate our two-dimensional RDMA scheduling. For the vertical dimension, we use the weighted min-max ratio (WMMR) $\frac{\min(x_i/w_i)}{\max(x_i/w_i)}$ [92] as our bandwidth fairness metric (the closer to 1, the better), where x_i is the bandwidth consumption of the application i , and w_i is its weight. We set the weight proportionally to the average bandwidth of each application when running individually. Our vertical scheduling achieves an overall of **0.88** WMMR.

The horizontal dimension (*i.e.*, priority scheduling with timeliness) is our focus here because interference between prefetching and demand swapping is a unique challenge we overcome in this work. We ran Spark GraphX-CC with the three native applications. Figure 13 compares the latency of sync vs. async swap-in requests with and without the horizontal scheduling of RDMA.

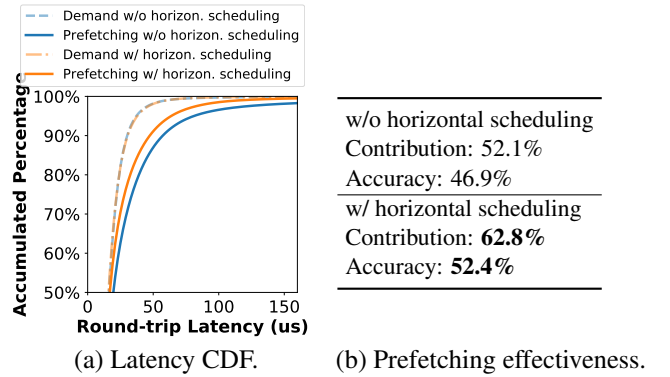


Figure 13: Horizontal scheduling effectiveness for GraphX-CC: (a) prefetching latency reduced, and (b) prefetching contribution and accuracy improved.

As shown, our scheduler does *not* incur overhead for the synchronous, demand requests but reduces the (90th percentile) latency of the asynchronous prefetching requests by $\sim 5\%$. Note that these results were obtained with Canvas’s two-tier prefetcher enabled, which already generates precise prefetching requests. With the LEAP prefetcher, the (90th percentile) latency reduction can be as high as $9\times$. To understand how the latency reduction improves prefetching effectiveness, we have also compared the prefetching contribution and accuracy with and without the horizontal scheduling, as shown in Figure 13(b). Due to the high timeliness requirement of prefetching requests, even 5% latency reduction can lead to noticeable improvements in prefetching—*e.g.*, the contribution/accuracy of GraphX-CC increases by **10.7%** and **5.5%** on top of the two-tier prefetcher—which ultimately translate to a **7-12%** overall improvement.

6.4 Other Analysis

We co-ran these applications with adaptive optimizations enabled but no isolation. In this case, we observed large per-

formance variations again. For example, when co-running Memcached, Snappy and XGBoost respectively with Spark-LR and Spark-PR, XGBoost is 15% and 2% slower than co-running with isolation. On the other hand, we ran XGBoost *individually* with all adaptive optimizations; this led to an 8% improvement over XGBoost co-run with other applications on Canvas. The win comes primarily from the exclusive use of the (unsaturated) RDMA bandwidth.

7 Related Work

Remote Memory. The past few years have seen a proliferation of remote-memory systems that built on the kernel’s swap mechanisms (including recent works such as LegoOS [87], Infiniswap [37], Fastswap [8], and Semeru [100] as well as earlier attempts [30, 6, 29, 32, 26, 43, 58]). Remote memory is part of a general trend of resource disaggregation in datacenters [41, 19, 34, 13, 10, 63, 62, 55, 7, 9, 79, 91], which holds the promise of improving resource utilization and simplifying new hardware adoption. Under disaggregated memory, application data are all stored on memory servers, making swap interference a more serious problem.

Resource Isolation. Interference exists in a wide variety of settings [27, 66, 106] and resource isolation is crucial for delivering reliable performance for user workloads. There is a large body of work on isolation of various kinds of resources including compute time [61, 14, 23], processor caches [33, 54, 101], memory bandwidth [64, 65, 68, 47, 102], I/O bandwidth [38, 92, 67, 71, 93, 99, 103], network bandwidth [12, 39, 35, 90, 83, 74, 50], congestion control [25, 42], as well as CPU involved in network processing [56]. Techniques such as IX [15] and MTCP [49] isolate data-plane and application processing at the core granularity.

Prefetching. Prefetching is an extensively-studied idea, which has been used in the design of hardware cache [97, 40, 107, 96, 73], compilers [94, 60, 85, 82, 57, 31], as well as operating systems [98, 70]. Detecting spatial patterns [72] is a common way to prefetch data. For example, various hardware techniques [89, 51, 48] have been developed to identify patterns (*i.e.*, sequential or stride) in addresses accessed. LEAP [70] is a kernel prefetcher designed specifically for applications using remote memory. Swap interference can reduce the effectiveness of any existing prefetchers, let alone that none of them consider complex (semantic) patterns.

Early work such as [81, 18] proposes application-level prefetching for efficient file operations on slow disks. Our prefetcher is, however, designed for a new setting where applications trigger page faults frequently and read pages from fast remote memory, and have much tighter latency budgets. Furthermore, in our setting, “application” is the JVM (the managed runtime), and our prefetching is driven by semantics automatically inferred by the language runtime, without needing any hints/annotations from application developers.

RDMA Optimizations. As RDMA gains popularity, there is a body of recent work on RDMA scheduling [84, 88]

and scalability improvement [95, 22, 53, 52]. Existing techniques focus on solving scalability problems when RDMA NICs are shared among multiple clients. Canvas extends existing fair scheduling and priority-based scheduling techniques to schedule prefetching/demand requests.

8 Conclusion

Canvas is a redesigned swap system that provides strong resource isolation and adaptive optimizations. It is also built with a set of optimizations that dynamically adapt their strategies to each application’s behavior/need based on the isolated swap paths. Canvas significantly improves the overall performance and minimizes swap-induced performance variation.

References

- [1] Large Text Compression Benchmark.
- [2] NVMe over fabrics. <http://community.mellanox.com/s/article/what-is-nvme-over-fabrics-x>.
- [3] Libsvm data: Classification. <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>, 2012.
- [4] Memcached - a distributed memory object caching system. <http://memcached.org>, 2020.
- [5] Wikipedia networks data. <http://konect.uni-koblenz.de/networks/>, 2020.
- [6] M. K. Aguilera, N. Amit, I. Calciu, X. Deguillard, J. Gandhi, S. Novakovic, A. Ramanathan, P. Subrahmanyam, L. Suresh, K. Tati, R. Venkatasubramanian, and M. Wei. Remote regions: A simple abstraction for remote memory. In *USENIX ATC*, pages 775–787, 2018.
- [7] M. K. Aguilera, K. Keeton, S. Novakovic, and S. Singhal. Designing far memory data structures: Think outside the box. In *HotOS*, pages 120–126, 2019.
- [8] E. Amaro, C. Branner-Augmon, Z. Luo, A. Ousterhout, M. K. Aguilera, A. Panda, S. Ratnasamy, and S. Shenker. Can far memory improve job throughput? In *EuroSys*, 2020.
- [9] S. Angel, M. Nanavati, and S. Sen. Disaggregation and the application. In *HotCloud*, 2020.
- [10] K. Asanović, R. Bodik, B. C. Catanzaro, J. J. Gebis, P. Husbands, K. Keutzer, D. A. Patterson, W. L. Plishker, J. Shalf, S. W. Williams, and K. A. Yelick. The landscape of parallel computing research: A

- view from berkeley. Technical Report UCB/EECS-2006-183, EECS Department, University of California, Berkeley, Dec 2006.
- [11] P. Baldi, P. Sadowski, and D. Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature communications*, 5(1):1–9, 2014.
 - [12] H. Ballani, P. Costa, T. Karagiannis, and A. Rowstron. Towards predictable datacenter networks. In *SIGCOMM*, pages 242–253, 2011.
 - [13] L. A. Barroso. Warehouse-scale computing: Entering the teenage decade. In *ISCA*, 2011.
 - [14] D. B. Bartolini, F. Sironi, D. Sciuto, and M. D. Santambrogio. Automated fine-grained cpu provisioning for virtual machines. *ACM Trans. Archit. Code Optim.*, 11(3), July 2014.
 - [15] A. Belay, G. Prekas, A. Klimovic, S. Grossman, C. Kozyrakis, and E. Bugnion. IX: A protected dataplane operating system for high throughput and low latency. In *OSDI*, pages 49–65, 2014.
 - [16] D. Bertsekas and R. Gallager. *Data Networks (2nd Ed.)*. Prentice-Hall, Inc., USA, 1992.
 - [17] I. Calciu, M. T. Imran, I. Puddu, S. Kashyap, H. A. Maruf, O. Mutlu, and A. Kolli. Rethinking software runtimes for disaggregated memory. In *ASPLOS*, pages 79–92, 2021.
 - [18] P. Cao, E. W. Felten, A. R. Karlin, and K. Li. Implementation and performance of integrated application-controlled file caching, prefetching, and disk scheduling. *ACM Trans. Comput. Syst.*, 14(4):311–343, Nov. 1996.
 - [19] A. Carbonari and I. Beschastnikh. Tolerating faults in disaggregated datacenters. In *HotNets-XVI*, pages 164–170, 2017.
 - [20] T. Chen and C. Guestrin. XGBoost: A scalable tree boosting system. In *KDD*, pages 785–794, 2016.
 - [21] T. Chen and C. Guestrin. extreme gradient boosting for applied machine learning. <https://xgboost.readthedocs.io/en/latest/>, 2021.
 - [22] Y. Chen, Y. Lu, and J. Shu. Scalable RDMA RPC on reliable connection with efficient resource sharing. In *EuroSys*, 2019.
 - [23] L. Cherkasova, D. Gupta, and A. Vahdat. Comparison of the three cpu schedulers in xen. *SIGMETRICS Perform. Eval. Rev.*, 35(2):42–51, 2007.
 - [24] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, and R. Sears. Benchmarking cloud serving systems with ycsb. In *Proceedings of the 1st ACM Symposium on Cloud Computing, SoCC '10*, page 143–154, New York, NY, USA, 2010. Association for Computing Machinery.
 - [25] B. Cronkite-Ratcliff, A. Bergman, S. Vargaftik, M. Ravi, N. McKeown, I. Abraham, and I. Keslassy. Virtualized congestion control. In *SIGCOMM*, pages 230–243, 2016.
 - [26] M. D. Dahlin, R. Y. Wang, T. E. Anderson, and D. A. Patterson. Cooperative caching: Using remote client memory to improve file system performance. In *OSDI*, 1994.
 - [27] C. Delimitrou and C. Kozyrakis. Bolt: I know what you did last summer... in the cloud. In *ASPLOS*, pages 599–613, 2017.
 - [28] A. Demers, S. Keshav, and S. Shenker. Analysis and simulation of a fair queueing algorithm. *SIGCOMM Comput. Commun. Rev.*, 19(4):1–12, Aug. 1989.
 - [29] M. J. Feeley, W. E. Morgan, E. P. Pighin, A. R. Karlin, H. M. Levy, and C. A. Thekkath. Implementing global memory management in a workstation cluster. In *SOSP*, pages 201–212, 1995.
 - [30] E. Felten and J. Zahorjan. Issues in the implementation of a remote memory paging system. In *University of Washington CSE TR CSE TR*, 1991.
 - [31] M. Ferdman, C. Kaynak, and B. Falsafi. Proactive instruction fetch. In *MICRO*, pages 152–162, 2011.
 - [32] M. D. Flouris and E. P. Markatos. The network ramdisk: Using remote memory on heterogeneous nodes. *Cluster Computing*, 2(4), Dec 1999.
 - [33] L. Funaro, O. A. Ben-Yehuda, and A. Schuster. Ginseng: Market-driven LLC allocation. In *USENIX ATC*, pages 295–308, 2016.
 - [34] P. X. Gao, A. Narayan, S. Karandikar, J. Carreira, S. Han, R. Agarwal, S. Ratnasamy, and S. Shenker. Network requirements for resource disaggregation. In *OSDI*, pages 249–264, 2016.
 - [35] A. Ghodsi, M. Zaharia, B. Hindman, A. Konwinski, S. Shenker, and I. Stoica. Dominant resource fairness: Fair allocation of multiple resource types. In *NSDI*, pages 323–336, 2011.
 - [36] Google. Google’s fast compressor/decompressor. <https://github.com/google/snappy>, 2020.

- [37] J. Gu, Y. Lee, Y. Zhang, M. Chowdhury, and K. G. Shin. Efficient memory disaggregation with infiniswap. In *NSDI*, pages 649–667, 2017.
- [38] A. Gulati, A. Merchant, and P. J. Varman. MClock: Handling throughput variability for hypervisor IO scheduling. In *OSDI*, pages 437–450, 2010.
- [39] C. Guo, G. Lu, H. J. Wang, S. Yang, C. Kong, P. Sun, W. Wu, and Y. Zhang. SecondNet: A data center network virtualization architecture with bandwidth guarantees. In *Co-NEXT*, 2010.
- [40] Y. Guo. *Compiler-Assisted Hardware-Based Data Prefetching for next Generation Processors*. PhD thesis, 2007.
- [41] S. Han, N. Egi, A. Panda, S. Ratnasamy, G. Shi, and S. Shenker. Network support for resource disaggregation in next-generation datacenters. In *HotNets*, pages 10:1–10:7, 2013.
- [42] K. He, E. Rozner, K. Agarwal, Y. J. Gu, W. Felter, J. Carter, and A. Akella. AC/DC TCP: Virtual congestion control enforcement for datacenter networks. In *SIGCOMM*, pages 244–257, 2016.
- [43] L. Iftode, K. Li, and K. Petersen. Memory servers for multicomputers. In *Digest of Papers. Compcon Spring*, pages 538–547, Feb 1993.
- [44] Intel. Batch allocation for swap entries. <https://github.com/torvalds/linux/commit/ed43af10975eef7e>, 2020.
- [45] Intel. Memcontrol: Charge swap-in pages to cgroup. <https://github.com/torvalds/linux/commit/4c6355b25e8bb83c>, 2020.
- [46] Intel. Per-core cluster allocation. <https://github.com/torvalds/linux/commit/490705888107c3ed>, 2020.
- [47] R. Iyer, L. Zhao, F. Guo, R. Illikkal, S. Makineni, D. Newell, Y. Solihin, L. Hsu, and S. Reinhardt. QoS policies and architecture for cache/memory in CMP platforms. In *SIGMETRICS*, pages 25–36, 2007.
- [48] A. Jain and C. Lin. Linearizing irregular memory accesses for improved correlated prefetching. In *MICRO*, pages 247–259, 2013.
- [49] E. Y. Jeong, S. Woo, M. Jamshed, H. Jeong, S. Ihm, D. Han, and K. Park. MTCP: A highly scalable user-level TCP stack for multicore systems. In *NSDI*, pages 489–502, 2014.
- [50] V. Jeyakumar, M. Alizadeh, D. Mazières, B. Prabhakar, and C. Kim. EyeQ: Practical network performance isolation for the multi-tenant cloud. In *HotCloud*, 2012.
- [51] D. Joseph and D. Grunwald. Prefetching using markov predictors. In *ISCA*, pages 252–263, 1997.
- [52] A. Kalia, M. Kaminsky, and D. G. Andersen. Design guidelines for high performance RDMA systems. In *USENIX ATC*, pages 437–450, 2016.
- [53] A. Kalia, M. Kaminsky, and D. G. Andersen. FaSST: Fast, scalable and simple distributed transactions with two-sided (RDMA) datagram RPCs. In *OSDI*, pages 185–201, 2016.
- [54] H. Kasture and D. Sanchez. Ubik: Efficient cache sharing with strict qos for latency-critical workloads. In *ASPLOS*, pages 729–742, 2014.
- [55] K. Keeton. The Machine: An architecture for memory-centric computing. In *ROSS*, 2015.
- [56] J. Khalid, E. Rozner, W. Felter, C. Xu, K. Rajamani, A. Ferreira, and A. Akella. Iron: Isolating network-based cpu in container environments. In *NSDI*, pages 313–328, 2018.
- [57] A. Kolli, A. Saidi, and T. F. Wenisch. RDIP: Return-address-stack directed instruction prefetching. In *MICRO*, pages 260–271, 2013.
- [58] S. Koussih, A. Acharya, and S. Setia. Dodo: a user-level system for exploiting idle memory in workstation clusters. In *HPDC*, pages 301–308, Aug 1999.
- [59] A. Lagar-Cavilla, J. Ahn, S. Souhlal, N. Agarwal, R. Burny, S. Butt, J. Chang, A. Chaugule, N. Deng, J. Shahid, G. Thelen, K. A. Yurtsever, Y. Zhao, and P. Ranganathan. Software-defined far memory in warehouse-scale computers. In *ASPLOS*, pages 317–330, 2019.
- [60] C. Lattner and V. Adve. Automatic pool allocation: improving performance by controlling data structure layout in the heap. In *PLDI*, pages 129–142, 2005.
- [61] T. Li, D. Baumberger, and S. Hahn. Efficient and scalable multiprocessor fair scheduling using distributed weighted round-robin. In *PPoPP*, pages 65–74, 2009.
- [62] K. Lim, J. Chang, T. Mudge, P. Ranganathan, S. K. Reinhardt, and T. F. Wenisch. Disaggregated memory for expansion and sharing in blade servers. In *ISCA*, pages 267–278, 2009.

- [63] K. Lim, Y. Turner, J. R. Santos, A. AuYoung, J. Chang, P. Ranganathan, and T. F. Wenisch. System-level implications of disaggregated memory. In *HPCA*, pages 1–12, 2012.
- [64] L. Liu, Z. Cui, M. Xing, Y. Bao, M. Chen, and C. Wu. A software memory partition approach for eliminating bank-level interference in multicore systems. In *PACT*, pages 367–376, 2012.
- [65] L. Liu, Y. Li, Z. Cui, Y. Bao, M. Chen, and C. Wu. Going vertical in memory management: Handling multiplicity by multi-policy. In *ISCA*, pages 169–180, 2014.
- [66] D. Lo, L. Cheng, R. Govindaraju, P. Ranganathan, and C. Kozyrakis. Improving resource efficiency at scale with heracles. *ACM Trans. Comput. Syst.*, 34(2), 2016.
- [67] L. Lu, Y. Zhang, T. Do, S. Al-Kiswani, A. C. Arpaci-Dusseau, and R. H. Arpaci-Dusseau. Physical disentanglement in a container-based file system. In *OSDI*, pages 81–96, 2014.
- [68] J. Ma, X. Sui, N. Sun, Y. Li, Z. Yu, B. Huang, T. Xu, Z. Yao, Y. Chen, H. Wang, L. Zhang, and Y. Bao. Supporting differentiated services in computers via programmable architecture for resourcing-on-demand (PARD). In *ASPLOS*, pages 131–143, 2015.
- [69] M. Maas, K. Asanović, T. Harris, and J. Kubiawicz. Taurus: A holistic language runtime system for coordinating distributed managed-language applications. In *ASPLOS*, pages 457–471, 2016.
- [70] H. A. Maruf and M. Chowdhury. Effectively prefetching remote memory with Leap. In *USENIX ATC*, pages 843–857, 2020.
- [71] J. C. McCullough, J. Dunagan, A. Wolman, and A. C. Snoeren. Stout: An adaptive interface to scalable cloud storage. In *USENIX ATC*, 2010.
- [72] M. K. McKusick, W. N. Joy, S. J. Leffler, and R. S. Fabry. A fast file system for UNIX. *ACM Trans. Comput. Syst.*, 2(3):181–197, 1984.
- [73] S. Mittal. A survey of recent prefetching techniques for processor caches. *ACM Comput. Surv.*, 49(2), 2016.
- [74] Y. Mundada, A. Ramachandran, and N. Feamster. Silverline: Data and network isolation for cloud services. In *HotCloud*, 2011.
- [75] C. Navasca, C. Cai, K. Nguyen, B. Demsky, S. Lu, M. Kim, and G. H. Xu. Gerenuk: Thin computation over big native data using speculative program transformation. In *SOSP*, pages 538–553, 2019.
- [76] K. Nguyen, L. Fang, C. Navasca, G. Xu, B. Demsky, and S. Lu. Skyway: Connecting managed heaps in distributed big data systems. In *ASPLOS*, pages 56–69, 2018.
- [77] K. Nguyen, L. Fang, G. Xu, B. Demsky, S. Lu, S. Alamian, and O. Mutlu. Yak: A high-performance big-data-friendly garbage collector. In *OSDI*, pages 349–365, 2016.
- [78] K. Nguyen, K. Wang, Y. Bu, L. Fang, J. Hu, and G. Xu. Facade: A compiler and runtime for (almost) object-bounded big data applications. In *ASPLOS*, pages 675–690, 2015.
- [79] A. Ousterhout, J. Fried, J. Behrens, A. Belay, and H. Balakrishnan. Shenango: Achieving high CPU efficiency for latency-sensitive datacenter workloads. In *NSDI*, pages 361–378, 2019.
- [80] A. Parekh and R. Gallager. A generalized processor sharing approach to flow control in integrated services networks: the single-node case. *IEEE/ACM Transactions on Networking*, 1(3):344–357, 1993.
- [81] R. H. Patterson, G. A. Gibson, E. Ginting, D. Stodolsky, and J. Zelenka. Informed prefetching and caching. In *SOSP*, pages 79–95, 1995.
- [82] L. Peled, S. Mannor, U. Weiser, and Y. Etsion. Semantic locality and context-based prefetching using reinforcement learning. In *ISCA*, pages 285–297, 2015.
- [83] L. Popa, A. Krishnamurthy, S. Ratnasamy, and I. Stoica. Faircloud: Sharing the network in cloud computing. In *SIGCOMM*, pages 187–198, 2012.
- [84] H. Qiu, X. Wang, T. Jin, Z. Qian, B. Ye, B. Tang, W. Li, and S. Lu. Toward effective and fair RDMA resource sharing. In *APNet*, pages 8–14, 2018.
- [85] R. M. Rabbah, H. Sandanagobalane, M. Ekpinyapong, and W.-F. Wong. Compiler orchestrated prefetching via speculation and predication. In *ASPLOS*, pages 189–198, 2004.
- [86] Z. Ruan, M. Schwarzkopf, M. K. Aguilera, and A. Belay. AIFM: High-performance, application-integrated far memory. In *OSDI*, pages 315–332, 2020.
- [87] Y. Shan, Y. Huang, Y. Chen, and Y. Zhang. LegoOS: A disseminated, distributed OS for hardware resource disaggregation. In *OSDI*, pages 69–87, 2018.

- [88] D. Shen, J. Luo, F. Dong, X. Guo, K. Wang, and J. C. S. Lui. Distributed and optimal rdma resource scheduling in shared data center networks. In *INFO-COM*, pages 606–615, 2020.
- [89] T. Sherwood, S. Sair, and B. Calder. Predictor-directed stream buffers. In *MICRO*, pages 42–53, 2000.
- [90] A. Shieh, S. Kandula, A. Greenberg, and C. Kim. Seawall: Performance isolation for cloud datacenter networks. In *HotCloud*, 2010.
- [91] V. Shrivastav, A. Valadarsky, H. Ballani, P. Costa, K. S. Lee, H. Wang, R. Agarwal, and H. Weather- spoon. Shoal: A network architecture for disaggregated racks. In *NSDI*, pages 255–270, 2019.
- [92] D. Shue, M. J. Freedman, and A. Shaikh. Performance isolation and fairness for multi-tenant cloud storage. In *OSDI*, pages 349–362, 2012.
- [93] E. Thereska, H. Ballani, G. O’Shea, T. Karagian- nis, A. Rowstron, T. Talpey, R. Black, and T. Zhu. IOFlow: A software-defined storage architecture. In *SOSP*, pages 182–196, 2013.
- [94] Tien-Fu Chen and Jean-Loup Baer. Effective hardware-based data prefetching for high- performance processors. *IEEE Transactions on Com- puters*, 44(5):609–623, 1995.
- [95] S.-Y. Tsai and Y. Zhang. LITE kernel RDMA support for datacenter applications. In *SOSP*, pages 306–324, 2017.
- [96] S. P. Vander Wiel and D. J. Lilja. When caches aren’t enough: data prefetching techniques. *Com- puter*, 30(7):23–30, 1997.
- [97] S. P. Vander Wiel and D. J. Lilja. A compiler-assisted data prefetch controller. In *Proceedings 1999 IEEE International Conference on Computer Design: VLSI in Computers and Processors*, pages 372–377, 1999.
- [98] G. M. Voelker, E. J. Anderson, T. Kimbrel, M. J. Fee- ley, J. S. Chase, A. R. Karlin, and H. M. Levy. Im- plementing cooperative prefetching and caching in a globally-managed memory system. In *SIGMETRICS*, pages 33–43, 1998.
- [99] M. Wachs, M. Abd-El-Malek, E. Thereska, and G. R. Ganger. Argon: Performance insulation for shared storage servers. In *FAST*, 2007.
- [100] C. Wang, H. Ma, S. Liu, Y. Li, Z. Ruan, K. Nguyen, M. D. Bond, R. Netravali, M. Kim, and G. H. Xu. Se- meru: A memory-disaggregated managed runtime. In *14th USENIX Symposium on Operating Systems De- sign and Implementation (OSDI 20)*, pages 261–280. USENIX Association, Nov. 2020.
- [101] X. Wang and J. F. Martínez. ReBudget: Trading off efficiency vs. fairness in market-based multicore re- source allocation via runtime budget reassignment. In *ASPLOS*, pages 19–32, 2016.
- [102] H. Yang, A. Breslow, J. Mars, and L. Tang. Bubble- Flux: Precise online qos management for increased utilization in warehouse scale computers. In *ISCA*, pages 607–618, 2013.
- [103] S. Yang, T. Harter, N. Agrawal, S. S. Kowsalya, A. Kr- ishnamurthy, S. Al-Kiswany, R. T. Kaushik, A. C. Arpaci-Dusseau, and R. H. Arpaci-Dusseau. Split- level i/o scheduling. In *SOSP*, pages 474–489, 2015.
- [104] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica. Spark: Cluster computing with working sets. *HotCloud*, page 10, Berkeley, CA, USA, 2010.
- [105] L. Zhang. A new architecture for packet switch- ing network protocols. Technical report, MAS- SACHUSETTS INST OF TECH CAMBRIDGE LAB FOR COMPUTER SCIENCE, 1989.
- [106] W. Zhang, S. Rajasekaran, S. Duan, T. Wood, and M. Zhuy. Minimizing interference and maximizing progress for hadoop virtual machines. *SIGMETRICS Perform. Eval. Rev.*, 42(4):62–71, 2015.
- [107] D. F. Zucker, R. B. Lee, and M. J. Flynn. Hardware and software cache prefetching techniques for MPEG benchmarks. *IEEE Transactions on Circuits and Sys- tems for Video Technology*, 10(5):782–796, 2000.