

End the Senseless Killing: Improving Memory Management for Mobile Operating Systems

Niel Lebeck
University of Washington

Arvind Krishnamurthy
University of Washington

Henry M. Levy
University of Washington

Irene Zhang
Microsoft Research

Abstract

To ensure low-latency memory allocation, mobile operating systems kill applications instead of swapping memory to disk. This design choice shifts the burden of managing over-utilized memory to application programmers, requiring them to constantly checkpoint their application state to disk. This paper presents Marvin, a new memory manager for mobile platforms that efficiently supports swapping while meeting the strict performance requirements of mobile apps. Marvin’s swap-enabled language runtime is co-designed with OS-level memory management to avoid common pitfalls of traditional swap mechanisms. Its key features are: (1) a new swap mechanism, called *ahead-of-time (AOT) swap*, which pre-writes memory to disk, then harvests it quickly when needed, (2) a modified bookmarking garbage collector that avoids swapping in unused memory, and (3) an object-granularity working set estimator. Our experiments show that Marvin can run more than 2x as many concurrent apps as Android, and that Marvin can reclaim memory over 60x faster than Android with a Linux swap file can allocate memory under memory pressure.

1 Introduction

Over the past decade, mobile apps have become bigger and more complex [28], far outpacing increases in mobile device memory [4]. This trend has increased memory pressure on mobile operating systems as apps compete for limited space. Going forward, mobile OSes must more efficiently share memory across demanding apps, or user experience will suffer.

Unfortunately, while mobile apps have become more sophisticated, mobile memory management remains in its infancy. Today’s popular mobile OSes set a fixed upper bound on memory for each running application (e.g., 1.4GB for iOS running on an iPhone X [36] and 512MB for Android running on a Google Pixel XL). They never overcommit memory; instead, they kill running applications and restart them later.

This simplistic approach worked well when mobile apps were small and largely stateless. However, it is unsustainable

as mobile OSes replace desktop ones (e.g., Android is now the most used OS in the world) and mobile apps replace desktop counterparts (e.g., Google Docs and Word Online replacing Microsoft Word). Today’s apps already do not fit into their memory allocation, so they manually swap objects between memory and local storage or use libraries to meet their needs [11]. Because apps are increasingly likely to be killed due to memory pressure, they must also continuously save execution state to disk and strive to minimize their start-up times to cope with frequent restarts. Despite significant engineering effort [21, 27, 32], it still takes several seconds to kill and restart popular apps.

Improving mobile memory management is difficult. Mobile apps run in high-level language runtimes (e.g., Swift, ART), which limit OS insight (e.g., working set estimation is impossible) and are notoriously difficult for OS-level memory managers to work with [20]. Further, mobile apps often allocate large amounts of memory quickly (e.g., when starting, or for cloud downloads); unless the OS keeps a large pool of free memory, this is easier to accommodate by killing entire applications. Finally, touch-based interfaces impose strict latency requirements, which swapping to disk cannot meet.

In this paper, we improve mobile memory management with a key observation: unlike other operating systems, mobile OSes run all of their apps in a *common* language runtime. For example, all apps running on Android must run in the Android Runtime (ART). This difference lets us co-design the language runtime to assist the mobile OS in optimizing memory management instead of hindering it. Due to its knowledge of memory usage, the language runtime becomes an ideal place for mechanisms that can better manage memory.

This paper demonstrates the value of leveraging the runtime for OS tasks. We present *Marvin*, a new memory manager for Android that efficiently supports memory overcommit. Marvin implements most memory management in the language runtime, which has more insight into an application’s memory usage. Marvin relies on the operating system only for cross-application resource allocation.

By integrating with the language runtime, Marvin can offer

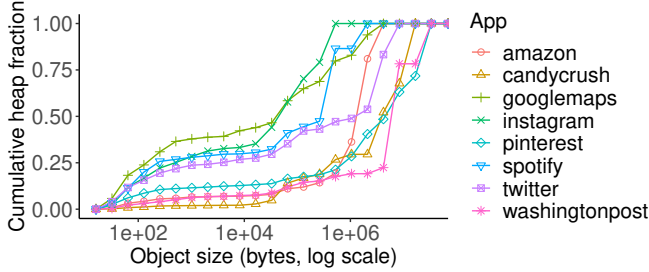


Figure 1: CDF of object size and heap percentage occupied by objects that size or smaller.

three new features that enhance memory management:

- A new swap mechanism, which performs *ahead-of-time* swapping to disk to avoid synchronous disk writes when reclaiming memory, and which leaves checkpointed objects in memory (unlike Linux’s kswapd [18]).
- A new object-level working set estimator, which separates garbage collector (GC) and app accesses, and avoids false sharing with object-level access tracking.
- A new bookmarking garbage collector [20], which tracks exact liveness data without accessing swapped-out objects.

We implement a prototype of Marvin by modifying the interpreter and compiler of the Android Runtime (ART). Experiments show that our Marvin prototype is able to run more than 2x as many concurrent apps as Android, and that Marvin can reclaim memory over 60x faster than Android with a Linux swap file can allocate memory under memory pressure.

2 Limitations of Modern Mobile OS Memory Resource Management

Although mobile OSes may be based on traditional OSes (e.g., Android and Linux), they diverge in two important ways: (1) for each app, they bound memory usage to a fraction of physical memory (e.g., 512MB on a 4GB device), rather than letting apps allocate as much memory as they need, and (2) they kill applications when physical memory runs out rather than overcommitting memory through paging or other mechanisms. To motivate our work, we ran experiments with popular apps that show the reasoning and cost for these design decisions. We ran all experiments on a Pixel XL phone with 4GB RAM and a quad-core Qualcomm Snapdragon 821 CPU.

2.1 Fixed Memory Allocation

Mobile OSes have poor insight into app memory usage. The runtime garbage collector regularly touches all objects and

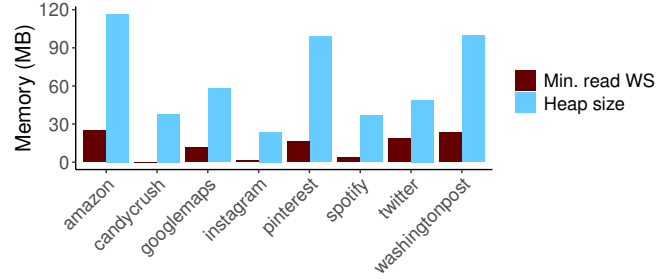


Figure 2: *The cost of fixed allocation.* Each bar shows the total Java heap size of a popular app alongside its minimum Java working set during active use.

moves objects for heap compaction, and the OS cannot distinguish this activity from app accesses. Mobile apps also access language-level objects, which vary in size, while the OS can only track memory accesses at page granularity. To understand the impact of object-level accesses on page-sized access tracking, we measured the size of objects in popular apps. Figure 1 shows a CDF of the size distribution. Most objects are not page-sized (e.g., up to 40% of objects are smaller than 4KB), so the OS cannot accurately track their usage.

Without good insight into app memory usage, today’s mobile OSes allocate all apps a fixed memory budget. On Android, this memory limit is the same whether an app is in the foreground or background. Android attempts to minimize the memory footprint of apps using techniques such as forking all apps from a single “zygote” process with copy-on-write pages. These techniques reduce duplication of framework data structures and shared libraries but do not impact application objects in the Java heap. Using Marvin’s object-level working set estimator, we measured the working set of popular apps. Figure 2 shows that although the heap footprint of these apps is large, their working sets actually account for a small fraction of their total heap size. This rarely accessed memory would be better utilized keeping other apps alive, rather than wasting space not being used.

Popular apps often have large memory footprints but small working set sizes because they cache as much as possible from the cloud. This caching improves performance, but it leads to poor memory utilization, and choosing the correct cache parameters is difficult [30]. Worse, modern applications frequently exceed their memory budgets. Coping with this problem requires apps to implement manual swap-to-storage, which adds significant programming complexity [10, 11]. While caching libraries like Glide [5] and Fresco [34] are helpful, they do not apply to all memory objects. Therefore, today’s apps use a complex combination of libraries and manually shuffling data between memory and disk.

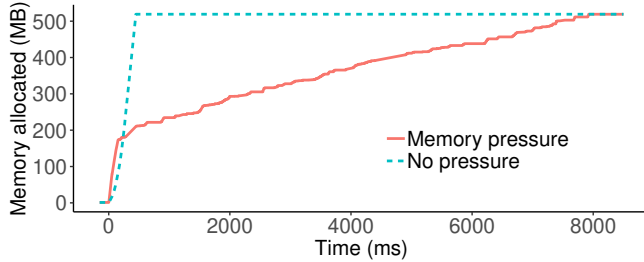


Figure 3: Progress over time of a memory allocation on Android with a swap file and memory pressure vs. no memory pressure.

2.2 No Memory Overcommit

Today’s mobile OSes kill applications rather than swapping to disk when physical memory runs out. They take this approach because mobile apps must respond to user input within hundreds of milliseconds, so traditional swap mechanisms, which place synchronous disk writes on the critical path, impose too much latency. To measure the effect of swapping on memory allocation, we enabled a Linux swap file on our Android test device [35], and we measured the amount of time required to allocate 512MB when the Android OS had free memory and when it had a swap file and memory pressure. Figure 3 shows the progress of the memory allocation over time. With memory available, the OS allocated all 512MB in 450ms; however, with a swap file and memory pressure, it took almost 8 seconds for the OS to allocate the same amount. Such high allocation latency would be unacceptable if an app were allocating memory in response to user input.

Unfortunately, killing and restarting apps comes at a cost. As shown in Figure 2, modern apps have large memory footprints, and a restarted app must fetch all of its cached data from the network or disk. We measured the amount of time needed to restart popular apps and compared it to that needed to fetch their entire checkpointed memory image from disk. As shown in Figure 4, restarting apps takes 4-27x longer than fetching the all of the app’s memory from disk.

The ability to kill and restart apps at any time also imposes a programming burden on app developers. Modern OSes give apps a limited time budget to perform cleanup before being killed. This limit leads apps to constantly write state to storage; in fact, Android encourages it [9]. Such constant checkpointing in response to app lifecycle events adds programming complexity, a challenge described in prior work [14]. Not only do app developers have to manage the checkpointing process, they have to correctly use a variety of mechanisms to do so with good performance [12]. The programming effort required to prepare for unexpected app deaths is an additional cost that app developers must pay.

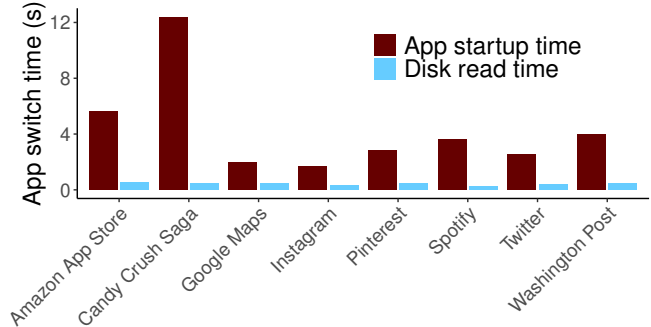


Figure 4: The cost of re-starting apps compared to reading their memory image from disk.

3 Our Approach

The primary barrier to improving memory resource management in mobile operating systems is the OS’s lack of insight into the language runtime. To overcome this barrier, we *co-designed* the language runtime and the mobile OS. Mobile operating systems are uniquely suited to such co-design because, unlike their desktop counterparts (e.g., Linux), they force all applications to use the same language runtime.

Marvin’s design focuses on Android (and the Android Runtime (ART)), which is now the most popular OS in the world [25]. Using the language runtime, Marvin manages memory entirely at object granularity, tracking, reclaiming and faulting in entire objects. This section describes in more detail the barriers to better memory resource management in a mobile OS and how Marvin addresses those challenges.

3.1 Object-Level Working Set Estimation

The first step to better memory resource management is a better understanding of each application’s working set. Thus, Marvin implements language-aware working set estimation in the language runtime, which tracks app reads and writes at object granularity. Marvin uses this mechanism to identify candidates for ahead-of-time swap (Section 3.2) and separate garbage collector accesses from app accesses (Section 3.3). Lacking hardware access bits to help with this tracking, Marvin implements software access tracking in both the ART interpreter and compiler, as modern mobile language runtimes run both interpreted and compiled code.

3.2 Ahead-of-time Swap

As noted, swapping to disk when the OS needs memory is not feasible for mobile OSes and their touch-based apps. Marvin takes a different approach. While traditional swapping mechanisms write to disk when memory is needed, Marvin uses a new *ahead-of-time* swap technique. This technique saves memory to disk *before* it is needed and then reclaims those

pages under memory pressure. Ahead-of-time swap separates swapping to disk from reclaiming memory; thus, we distinguish between *saved* objects, which have been copied to disk but still reside in memory, and *reclaimed* objects, which no longer reside in memory but are only on disk.

Swapping objects before they are needed leaves a large pool of clean memory that the OS can quickly reclaim and reallocate. While this technique lets apps continue using memory until the OS reclaims it, whenever the app dirties a page, the swap mechanism must update the on-disk copy before the OS can reclaim it. Due to this trade-off, Marvin prioritizes swapping objects that are infrequently or never written.

3.3 Bookmarking Garbage Collector

Like traditional swapping, ahead-of-time swapping is affected by friction with the language-level garbage collector. As noted by Hertz et al. [20], the garbage collector can inadvertently page in memory when walking the object heap to look for unused objects. With ahead-of-time swapping, the garbage collector can also inadvertently dirty pages when updating references, causing unnecessary writes to disk. Marvin solves this problem by integrating a modified bookmarking garbage collector [20] into the Android Runtime.

Marvin’s swap mechanism leaves *stubs* – analogous to bookmarks – for each reclaimed object that detail the objects’ references to other objects. Using these stubs, its garbage collector can process a reclaimed object during a mark-and-sweep run without faulting in the entire swapped object. Marvin’s swap mechanism can further optimize swapping from disk by dropping dead objects without faulting them in.

4 Marvin Overview

Marvin is a new mobile memory manager that supports flexible memory allocation between apps and memory overcommit through swapping. Marvin includes components in the language runtime and OS, which are co-designed to provide better memory management. This section overviews both.

4.1 Design Goals

Marvin’s design meets the following goals:

1. **Fast memory allocation.** Marvin must allocate memory quickly on-demand, avoiding disk accesses on the critical path for memory allocation.
2. **High memory utilization.** Marvin must provide the illusion of unlimited memory, provided working sets do not exceed the size of physical memory.
3. **Minimal overhead.** Marvin must impose low runtime overhead and require no app code changes.

While the last two goals are common to all memory management systems, existing mobile platforms sacrifice high memory utilization for fast memory allocation. Marvin aims to achieve all of these goals.

4.2 Marvin System Model

Marvin assumes a systems environment that meets three requirements: (1) all apps are written in a single managed language (e.g., Java), (2) all apps run in a single managed language runtime (e.g., ART), and (3) the runtime performs garbage collection or some form of automatic memory management. Marvin’s design targets Android, which meets all of these requirements. Android also runs some libraries using native code; however, Marvin is not needed to manage their memory because they do not have the same issues with OS-level memory management as managed languages.

Marvin’s optimizations could apply to other operating systems as well (e.g., iOS). For example, Swift uses automatic reference counting as an alternative to garbage collection, so it would require a bookmarking reference counter that can track references without faulting in the entire object.

Marvin runs unmodified Android apps on ARM64-based devices. Android distributes apps in a bytecode format called DEX. ART runs DEX bytecode directly in an interpreter and also compiles DEX to native ARM64 instructions both at install time (ahead-of-time, or AOT, compilation) and at runtime (just-in-time, or JIT, compilation). Marvin modifies both the interpreter and compiler.

4.3 Marvin Architecture

Marvin has two key components: (1) the *Marvin Kernel* (MK), a modified Android/Linux kernel, and (2) the *Marvin Runtime* (MRT), a modified ART. Most memory management occurs in MRT; it performs working set estimation, ahead-of-time swapping, and bookmarking garbage collection. MK’s sole responsibility is to balance memory allocation among apps by deciding when and from which app to reclaim memory.

Marvin performs working set estimation and swapping at object granularity; however, there is no CPU support for object-level access bits and memory faults. As a result, Marvin implements software object access tracking and faulting in the MRT interpreter and compiler. The interpreter marks access bits and checks for swapped-out objects as it runs DEX bytecode; the compiler inserts that functionality as additional ARM64 instructions. Marvin reserves four bytes in each object header and uses them to store swapping metadata and access bits. Implementing these features in software imposes overhead, which we quantify in Section 8. Hardware improvements in future mobile devices could reduce this overhead.

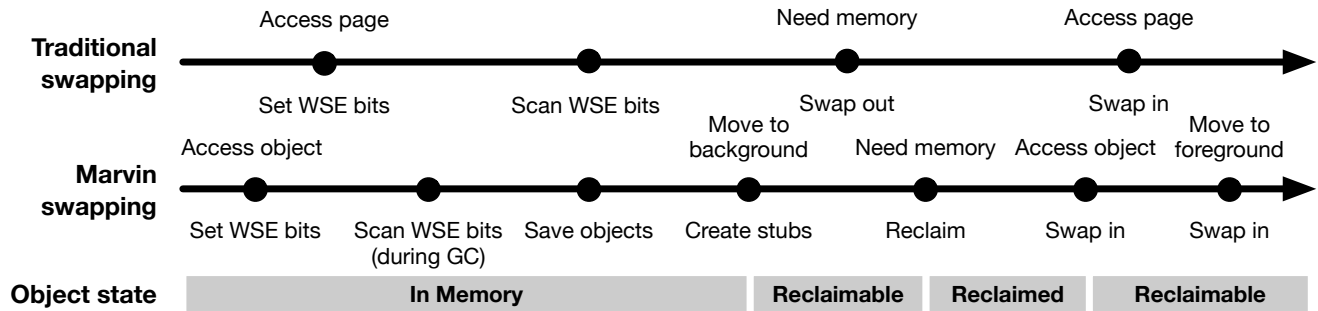


Figure 5: A timeline of actions performed by Marvin’s swap mechanism as compared to traditional (e.g., Linux) swap mechanisms. Events are listed above the timeline while Marvin’s actions in response are listed below.

4.4 Marvin Memory Management Timeline

Objects managed by Marvin move through several states over time, driven by app behavior and app lifecycle events. Figure 5 illustrates these events and states and compares Marvin’s swapping to a traditional swap mechanism. When an app first starts, MRT begins tracking its working set. It identifies objects that are suited for swapping by examining whether they are cold (have not been read or written recently by the app). MRT begins saving checkpoints of those objects to disk in the background. We refer to an object with a saved checkpoint as a *saved* object.

When the app moves from foreground to background, MRT pauses app threads and creates stubs, small proxy objects that add a layer of indirection over swap candidate objects. Stubs ensure that Marvin can intercept accesses to objects and fault them back in, if necessary. Once MRT creates a stub for an object, that object becomes *reclaimable*; the object is still memory-resident, but MK can reclaim its memory at any time. When MK reclaims an object, it enters the *reclaimed* state; only the object’s stub remains in memory, and the object’s checkpoint will need to be faulted back into memory before the object can be accessed again. The garbage collector uses only the stub and need not fault the object back into memory.

5 Marvin Core Mechanisms

As noted in Section 3, Marvin’s key features are ahead-of-time swap, language-aware working set estimation, and bookmarking garbage collection. Designing these features required addressing three challenges: adding a layer of indirection for object references, coordinating between the OS and runtime, and interposing on object accesses. This section describes Marvin’s mechanisms for addressing these challenges.

5.1 Stubs for Object Reference Indirection

According to the Java language specification, object references are opaque. However, in practice, object references in

the Android Runtime are direct pointers to the heap memory holding the referenced object. This design requires Marvin to inject a layer of indirection to implement features like software object faulting and the bookmarking garbage collector. Marvin creates this layer of indirection using special objects, called *stubs*. Each stub contains a pointer to its underlying object along with a copy of each reference held by the object. All references to an object point instead to its stub, and only the stub holds a pointer to its underlying object. Accessing an object through a stub adds overhead, so Marvin creates stubs only for objects that are cold and at least 2KB in size.

When creating a stub for an object, Marvin moves the object to a separate page-aligned region of memory and then redirects all references to the object to point to its stub instead. These tasks require that all app threads be paused. Therefore, they can be performed more efficiently if Marvin can create stubs for many objects at once, using a single scan of the heap to redirect all affected references. As a result, Marvin periodically executes a heap task that pauses all app threads and creates stubs, and it executes this heap task only when the app is in the background and its threads can be safely paused. Stubs are only created once for each object, so the cost is low, especially because Marvin does not restart apps frequently.

5.2 Reclamation Table for OS-Runtime Coordination

Modern mobile platforms have multicore processors that let system services run concurrently with apps. In this environment, the OS should be able to reclaim memory quickly from a running app without scheduling the app’s threads for execution. However, the OS cannot simply seize memory from an app whose threads are not scheduled—a pointer to the memory in question may be present in an app thread’s stack or registers, waiting to be used as soon as the thread is scheduled once again. As a result, the OS and runtime need a way to coordinate concurrent accesses to objects so the OS does not try to reclaim one that the runtime is accessing.

Marvin uses a shared-memory *reclamation table* to provide

this coordination. MRT populates the table with reclaimable objects, and MK uses it to identify memory to reclaim. Each reclamation table entry is a small, fixed-size data structure that holds the address of an object, its size, a set of flags indicating whether the object is memory-resident and the entry is valid, and a set of bits used for locking by the runtime and OS. To reclaim an object, MK first acquires an exclusive lock on the object's reclamation table entry. Similarly, whenever an app thread prepares to access an object, MRT acquires a shared lock on the reclamation table entry.

The reclamation table is necessary because it provides an agreed-upon location that MK can quickly scan to identify reclaimable objects. If the metadata in the reclamation table were instead stored inside the stubs themselves, then MK would need to scan MRT's entire Java heap to identify reclaimable objects, and stub headers would need to contain a magic number or provide some other way for the kernel to recognize them.

5.3 Object Access Interposition

All of Marvin's features require the runtime to interpose and perform specific tasks whenever an app accesses an object. On every object access, MRT must set read and write bits for working set estimation; check for the presence of a stub and redirect the object access through the stub if necessary; and fault in the object if it has been reclaimed. It must also set a dirty bit whenever an object is modified so the ahead-of-time swap mechanism knows which objects need to be saved, and it must update stubs whenever reference member variables in their corresponding objects change to support the bookmarking garbage collector.

Android apps execute both as DEX bytecode running in an interpreter and as compiled native code running directly on the hardware, and Marvin must interpose on all object accesses in both kinds of code. As a result, Marvin features a set of paired interpreter and compiler modifications that add the required object access interposition. For each additional task performed by the interpreter when it accesses an object, there is a corresponding change to the compiler, adding assembly instructions performing the same task to compiled code.

6 Marvin Memory Management

This section describes how we use Marvin's mechanisms (stubs, the reclamation table, and object access interposition) to design the features that make up Marvin's memory management system.

6.1 Working Set Estimation

MRT performs object-granularity working set estimation by maintaining two access bits in each object header, a read bit and a write bit, and scanning those access bits.

Setting access bits. MRT uses object access interposition to set an object's read and write bits whenever that object is read or written from either interpreted or compiled code. It avoids including garbage collector reads in its working set estimation by setting a flag in the object header when the garbage collector is visiting an object and leaving the read bit untouched if that flag is set.

Access tracking in MRT is performed on a best-effort basis to minimize its overhead: MRT uses non-atomic operations with relaxed memory ordering semantics when setting read and write bits. As a result, concurrent reads and writes to the same object could result in a lost update to one of the access bits. An update to the read bit could also be lost if an app thread reads an object that the garbage collector is processing. These optimizations may decrease swapping performance if the estimated and actual working sets differ significantly, but they do not affect correctness.

Scanning access bits. MRT periodically walks the heap and uses the Clock algorithm [8] to track each object's long-term usage. Each object header holds two four-bit *shift registers*, one for reads and the other for writes. The time between heap walks constitutes an access-tracking "round," and each shift register tracks whether the object was read or written in the last four rounds. During a heap walk, MRT updates an object's shift registers and then clears the object's access bits.

MRT piggybacks off of garbage collection to scan access bits, since GC requires walking the heap anyways. Some of ART's GCs only walk subsets of the heap, so MRT limits its access bit scanning to full-heap collections. It also periodically invokes GC to ensure up-to-date working set estimates in the absence of app activity.

Producing the working set. As MRT walks the heap and scans access bits, it tabulates the app's working set. Our current MRT implementation considers an object part of the working set if it has been written within the last four access-tracking rounds. The precise policy is an implementation detail that can be easily changed.

6.2 Ahead-of-Time Swapping

In Marvin's ahead-of-time swap mechanism, MK reclaims objects and decides which apps to target for reclamation. MRT performs all other functions, including saving object checkpoints to disk, restoring reclaimed objects, and preventing the operating system from reclaiming objects in use by app code.

Saving objects to disk. MRT identifies suitable objects for swapping (i.e., cold objects) using its working set estimation feature, and it proactively saves checkpoints of them to a swap file on disk so they can be reclaimed quickly under memory pressure. MRT saves objects to disk in a periodic heap task that runs on a background thread concurrently with app code.

After app code modifies an object, MRT must save an updated copy of that object to disk. It does not need to save the updated copy immediately as long as it prevents the kernel from reclaiming the object while it is “dirty.” To do so, MRT maintains a *dirty bit* in the object header. It uses object access interposition to set this dirty bit whenever app code writes to an object, and its object-saving heap task clears this dirty bit when saving the object to disk. MK checks dirty bits when looking for objects to reclaim and avoids reclaiming dirty objects. MRT and MK use strong memory-ordering semantics when reading and writing the dirty bit to ensure that no modifications to objects are lost.

MRT begins saving swap candidate objects to disk even before those objects have had stubs created for them. Once MRT creates a batch of stubs, those objects become immediately reclaimable without requiring further disk I/O.

Reclaiming objects. MK selects apps to target for reclamation and reclaims objects from the MRT instances corresponding to those apps. It never targets the foreground app, in order to avoid any swapping delays on the foreground app’s user interface (UI) thread. After selecting an MRT instance to target, MK scans the MRT instance’s reclamation table until it finds an entry for an object that is neither dirty nor locked by the runtime. MK then locks that entry and reclaims the object’s pages. It continues scanning the reclamation table and reclaiming objects until it has harvested the desired amount of memory from the MRT instance.

Each MRT instance ensures that MK does not reclaim an object currently being accessed by its app code. To do so, it uses object access interposition to detect whenever a reclaimable object is being read or written, and it locks the object’s reclamation table entry before the access and unlocks its entry after the access.

Restoring objects. MRT restores reclaimed objects either eagerly or on-demand. Either way, whenever MRT restores an object, it locks the object’s reclamation table entry, copies the saved checkpoint data of the object into memory, and copies any modified references from the object’s stub into the object itself. This last step is necessary because references in the stub may have been modified by the garbage collector while the object was not memory-resident.

Marvin’s eager object restoration uses app lifecycle information to restore objects before app code needs them. We implemented a simple eager restoration policy, where an MRT instance restores all reclaimed objects when it transitions to the foreground. This policy increases the transition delay in exchange for a guarantee that no swapping delays will block the foreground app’s UI thread. Our design is flexible and could support more advanced policies; for instance, the runtime could predict which objects are likely to be touched immediately after a foreground transition and restore those

objects first, trading off a shorter pause time in exchange for the risk of user-perceptible stuttering.

If an object has not been eagerly restored, MRT restores it on-demand when app code accesses it, a process that we call *software object faulting*. Whenever app code accesses a reclaimable object, MRT uses object access interposition to check if the object is memory-resident by inspecting a bit in its reclamation table entry. If not, MRT executes an object fault handler that calls into its C++ object restoration function.

6.3 Bookmarking Garbage Collector

A tracing garbage collector touches every object in the heap (or a subset of the heap), causing live objects to be swapped back into memory even if app code is not using them. Marvin’s garbage collector avoids touching reclaimed objects by storing an object’s references inside its stub and using the stub during the mark phase of garbage collection. Stubs play a similar role as bookmarks in the bookmarking collector [20].

During the mark phase, the garbage collector maintains a *mark stack* and repeatedly pops an object from the mark stack, marks all its references, and pushes those references onto the mark stack. Marvin’s garbage collector checks whether an object is a stub when it pops the object from the mark stack; if so, it reads the references off the stub instead of accessing the underlying object.

For the garbage collector to use stubs in place of their objects, MRT must ensure that the stub of a memory-resident object has up-to-date copies of the object’s references. It uses object access interposition to update the stub of a reclaimable object whenever Java code modifies one of the object’s reference member variables.

MRT must also properly clean up after any saved objects that are freed by the garbage collector. MRT records when saved objects have been freed by the garbage collector, and when the fraction of the swap file consisting of freed objects passes a set threshold (25% in our implementation), it compacts the swap file in a heap task. MRT also cleans up after reclaimable and reclaimed objects by checking whether an object being freed is a stub; if so, it deletes the reclamation table entry corresponding to the stub. If an object is reclaimable, MRT deletes the memory-resident copy of the underlying object; if the object is reclaimed, MRT simply marks its copy in the swap file for deletion without needing to fault it in.

6.4 Design Tradeoffs and Alternatives

By tracking working sets and faulting in objects in software at the runtime level, Marvin achieves a clean design, albeit with some drawbacks. First, Marvin cannot reclaim objects accessed by native libraries: native libraries have no way to detect stubs and no recourse for faulting in reclaimed objects. Second, software working set estimation and object faulting

add overhead, particularly to compiled code. We evaluate this overhead in Section 8.

Marvin moves almost all memory management into the runtime because we believe that the runtime’s better access to information about app behavior makes it better suited for managing memory. The functionality remaining in the kernel is the minimum required by existing Linux kernel design; if Marvin was built on top of an exokernel [13, 23], it could move even more functionality into the runtime. A variety of other designs are possible that split functionality between the runtime and kernel in different ways.

Kernel-level working set estimation would reduce the overhead of accessing objects, but it would suffer from false sharing if an app’s working set is mixed with unused objects across 4KB pages. Faulting in memory at the kernel level would similarly reduce object access overhead but would require more extensive re-design of the runtime garbage collector to avoid unnecessary swapping activity. By tying the granularity of memory management to the size of pages, kernel-level memory management will also become inflexible as large pages become more common and the disparity between object sizes and page sizes widens. In any case, kernel-level memory management would require some sort of ahead-of-time swap mechanism to satisfy the latency requirements of modern mobile platforms (Figure 3), and even adding ahead-of-time swap to the Linux kernel would require significant effort.

Marvin’s garbage collector is different from the original bookmarking collector [20] in that it maintains exact reachability information with stubs rather than conservatively storing approximate reachability information. The latter approach requires the garbage collector to perform less work when evicting pages and scanning the heap, but it can result in the heap being needlessly occupied with dead objects.

7 Marvin Prototype

We implemented a prototype of MRT by modifying ART on Android 7.1.1 (which includes Linux 3.18.31). Our implementation includes a modified version of ART’s ARM64 compiler, allowing our prototype to support Android devices with 64-bit ARM processors. Our changes to the ART codebase resulted in 3475 additional lines of code [38].

In addition to modifying ART, we made a small modification to the version of OpenJDK included with Android 7.1.1, namely, we added fields to the Object class definition to mirror the bytes added to the object header in ART. We also changed a source file in the Android framework (`ProcessList.java`) to increase a hard-coded limit on the number of concurrently running apps since Marvin is able to run more.

Our experiments require us to manually trigger reclamation, so we did not implement automated reclamation in the MK. However, our MRT implementation includes the reclamation table and performs all operations required to support kernel memory reclamation.

7.1 Object Access Interposition

We implemented MRT’s object access interposition by adding specialized functionality to the ART interpreter and compiler. This lets MRT interpose on object accesses from both DEX bytecode running in the interpreter and compiled OAT code running natively. The following section describes in detail how we modified each component.

MRT interpreter. The ART interpreter internally represents each Java object as a C++ *mirror object*, which it manipulates when executing DEX bytecode instructions that read or write an object. The mirror object’s type definition includes methods to read or write the data at a given offset within the object’s memory footprint, and the interpreter code calls these methods when executing DEX instructions. To add object access interposition to the interpreter, we modified the mirror object methods to implement Marvin’s features.

For example, to redirect object accesses through stubs and perform on-demand object faulting, we added a preamble macro to each mirror object method. The preamble first checks if the object is actually a stub. If so, it casts the `this` pointer to a stub, calls a method that locks the stub’s reclamation table entry (RTE), checks the RTE’s resident bit, and if the resident bit is cleared, calls a method to fault in the object from disk. The preamble then gets the address of the underlying object from the RTE and invokes the mirror object method on the underlying object. Finally, the preamble unlocks the RTE and returns the result of the mirror object method, if any.

ART contains multiple interpreter implementations, and the default is the “mterp interpreter,” an interpreter written in assembly. When the mterp interpreter executes DEX instructions that read or write an array, it directly accesses the array’s memory, bypassing the mirror object methods. To allow Marvin to interpose on array accesses, we instead use the “switch interpreter,” an interpreter written in C++ that calls the mirror object methods when executing array accesses.

MRT compiler. Java code in Android framework libraries and portions of app Java code execute as native code, which is compiled by the ART compiler either statically after installation or dynamically with just-in-time (JIT) compilation. We added object access interposition to this compiled code by modifying the compiler’s assembler to generate additional assembly instructions that implement Marvin’s features when it performs code generation for object accesses. We used ARM64 devices for testing and evaluation, so we added support for object access interposition to the ARM64 assembler.

Each operation described above for the interpreter’s implementation of stub redirection and object faulting has a corresponding block of ARM64 instructions in compiled code. The main difference is that when a stub is detected, the compiled code must explicitly overwrite the register holding the stub’s address with the address of the underlying “real object;” it

then loads the stub’s address back into that register when it is done with the object. In the common case, when an object is not a stub (or when it is, but its underlying “real object” is memory-resident), execution branches past many of the added object-faulting instructions.

7.2 Limitations and Potential Optimizations

Our MRT implementation is a research prototype and has some limitations as a result. One is instability when reclaiming objects from black-box commercial apps. When running apps that we create in Android Studio, MRT reliably and consistently reclaims and restores objects, but when running commercial apps with stub creation and reclamation enabled, it tends to crash. MRT’s object access interposition works correctly with commercial apps, so by disabling stub creation and reclamation, we can test its overhead and collect working set data. Our MRT prototype also does not support reclaiming objects with live JNI global references or directly accessing a reclaimable array’s memory through JNI.

Our implementation of object access interposition in the MRT compiler is unoptimized, and the per-object-access overhead of compiled code could be reduced with deeper compiler integration. We modified the ARM64 assembler, which translates intermediate representation (IR) instructions to ARM64 binary code. Our implementation generates ARM64 instructions performing object access interposition for every IR instruction that reads or writes an object, even though many of those instructions only need to execute when the object’s register allocation begins or ends. An optimized version of the compiler, with modifications at the IR level, could decrease both the execution and size overhead of compiled code.

The MRT compiler has other areas for optimization. For instance, we noticed situations where ARM64 parameter registers ($x0-x7$ and $d0-d7$) appear to be live but are not reported as such by the ART compiler’s `LocationSummary` class; therefore, we conservatively save all parameter registers to the stack when performing a procedure call. Saving only live parameter registers would further reduce code size overhead. In addition, the ARM64 assembler makes a maximum of two scratch registers available at any time, which required us to save backpointers and add more instructions to juggle required state among the limited available registers.

8 Evaluation

Our evaluation demonstrated that Marvin successfully met its design goals. It could: (1) quickly reclaim memory on-demand; (2) maintain high memory utilization by sharing memory among apps rather than killing them; and (3) achieve the previous two goals with low overhead and no app changes.

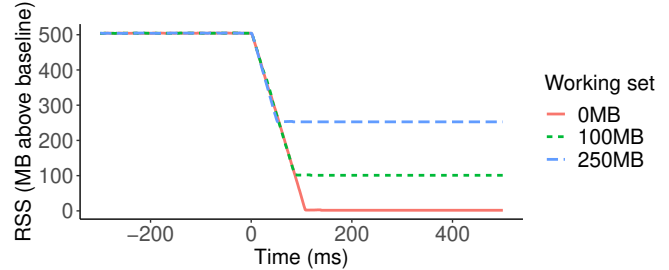


Figure 6: Memory usage as Marvin reclaims memory from a benchmark app with different working set sizes.

8.1 Evaluation Setup

We ran our experiments on a Google Pixel XL smartphone with 4GB RAM and a quad-core Qualcomm Snapdragon 821 CPU. The smartphone ran either the open-source release of Android 7.1.1 (AOSP tag `android-7.1.1_r57`) or our Marvin implementation based on that release. Both our Marvin implementation and our baseline Android build included a change to the Android framework to increase a hard-coded cap on the number of concurrently running apps.

Our experiments used a mix of synthetic apps for benchmarking and real-world apps. We built several synthetic workloads that simulate various memory footprints and working set sizes to measure their effect on Marvin compared to Android. We also used PCMark for Android, a commercial benchmark app based on real-world apps, to measure Marvin’s overhead on real apps [3]. We chose two benchmarks from the test suite (Writing 2.0 and Data Manipulation) since the remaining three test the performance of native libraries.

8.2 Memory Reclamation

Marvin must be able to quickly reclaim memory from running apps when a new or existing app needs to allocate large amounts of memory. Its ahead-of-time swap mechanism ensures that each MRT instance has a pool of clean memory that can be quickly reclaimed without swapping to disk. In this section, we measure the latency of reclaiming memory for apps with different working set sizes. Reclamation latency depends on the app’s working set size since Marvin can reclaim more memory from apps with smaller working set sizes.

For our prototype, we use `madvise` to return memory from MRT to the kernel. This design lets us trigger reclamation rather than waiting for memory pressure. Our MRT prototype reclaims memory when an app transitions to the background and then periodically while it is in the background.

Figure 6 shows memory usage over time for apps with a 500MB heap and differing working set sizes. RSS values shown are relative to the RSS reported for a minimal Android app with a single empty Activity (approx. 80MB). At time 0ms, MRT begins to return memory from the app to the OS.

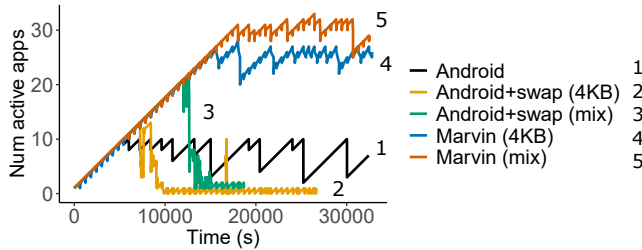


Figure 7: Count of active benchmark app instances over time. Marvin runs more than twice as many apps as regular Android before needing to kill any apps; on Android with a swap file, most apps are alive but inactive due to constant swapping activity.

Marvin returned 250MB of memory in 52ms and 500MB of memory in 108ms. In comparison, as shown in Figure 3, Android with a Linux swap file took nearly 8 seconds to free and allocate 500MB of memory under memory pressure. Using ahead-of-time swap let Marvin reclaim memory over 60x faster than Android with Linux swap could allocate the same amount of memory, allowing Marvin to meet the strict latency requirements of mobile apps.

8.3 Memory Utilization

To demonstrate Marvin’s more efficient memory manager, we ran multiple instances of an app with a large memory footprint and a limited working set, and we counted the number of *active* apps that were alive and making progress on their workloads. Each app had a 220MB heap filled with arrays, and it deleted and reallocated 20MB of those arrays every 5 seconds. We used two different heap compositions: one where the apps had heaps filled with 4KB arrays, and one where they had an even mix of 4KB and 1MB arrays (similar to the bimodal distribution of real apps in Figure 1). We consider an app “inactive” if it fails to perform a round of its workload for 20 seconds after the previous round; we consider it “active” once again if it succeeds in performing a round of its workload within 7 seconds of the previous round. We started a new app instance every 10 minutes to give Marvin time to perform background work. For unmodified Android, only the data for the apps with 4KB arrays is shown, because its behavior was nearly identical for the 4KB/1MB mix.

As shown in Figure 7, Marvin ran over 2x as many active apps concurrently as unmodified Android and over 1.5x-2x as Android with a Linux swap file enabled, where swapping left almost all apps unusable as the experiment went on. While baseline Android begins killing apps when physical memory runs out, Android with swap keeps more apps alive. However, without a bookmarking garbage collector, the system experienced constant swapping activity, which prevented most apps from making progress on their workloads. Our experimental

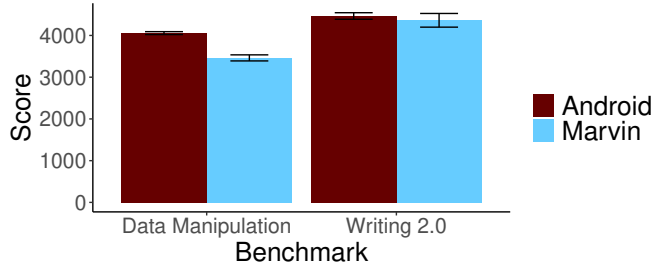


Figure 8: PCMark for Android benchmark results.

runs of Android with a swap file consistently ended early due to the device crashing.

Marvin made better use of device memory because it reclaimed unused memory from apps and used its bookmarking garbage collector to avoid touching that unused memory when running garbage collection. While Android only ran 10 apps concurrently, and Android with a swap file only briefly reached a maximum of 13 concurrent apps (4KB arrays) or 20 apps (4KB/1MB mix), Marvin ran 27 apps (4KB arrays) or 30 apps (4KB/1MB mix) concurrently. Marvin’s memory reclamation and bookmarking garbage collector let it execute 1.5-2x as many apps concurrently while neither killing apps nor suffering performance degradation.

8.4 Runtime Overhead

While Marvin provides better memory management, it comes with a set of trade-offs. This section quantifies Marvin’s four sources of overhead: (1) *execution time overhead* caused by Marvin’s object access interposition in compiled OAT code; (2) *increased compiled code size* due to object access interposition; (3) *CPU utilization overhead* caused by Marvin’s heap walks for working set estimation; and (4) *faulting overhead* when an app accesses a reclaimed object.

Execution time overhead of object access interposition. Native code produced by the MRT compiler has additional ARM64 instructions to support object access interposition. Some instructions (stub checks, dirty bit updates, and access-tracking bit updates) execute on every object access. Other instructions (indirecting object accesses through stubs and locking RTEs) execute only on accesses to reclaimable objects. We measured the overhead of the added instructions that apply to all object accesses using PCMark for Android, and we used a synthetic benchmark to illustrate the dependence of that overhead on the makeup of application code.

Figure 8 compares the performance of Marvin and unmodified Android on the PCMark benchmarks in our test set. Each bar shows the mean and standard deviation of five runs. We turned off stub creation and swapping when running PCMark, using the benchmarks to measure the overhead of the instructions added to every object access for stub checks and working

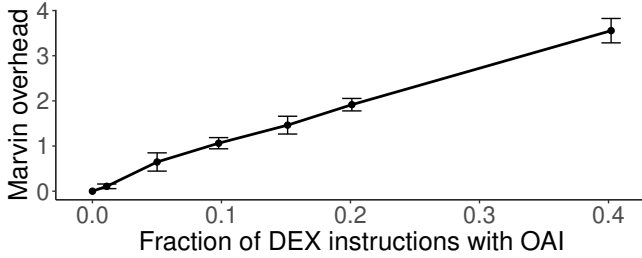


Figure 9: Overhead of Marvin for a synthetic workload with different proportions of object access interposition (OAI). The point (0,0) represents the theoretical scenario of running without any OAI, while other points show experimental results.

set estimation. Marvin’s score on the Writing 2.0 benchmark was nearly identical to Android’s, and its score on the Data Manipulation benchmark was only 15% lower. These scores show that Marvin’s overhead for accessing regular (i.e., non-reclaimable) objects is low for real-world apps.

Figure 9 explores the dependence of Marvin’s overhead on the DEX instruction mix of application code. The graph shows Marvin’s overhead executing a synthetic workload that performed a tunable proportion of object accesses (array reads and writes) and integer operations (addition and multiplication). Each point represents the mean and standard deviation of Marvin’s execution time overhead relative to Android for 40 iterations of the workload. For large proportions of object accesses, Marvin had relatively high overhead (e.g., 350% overhead for 40% object accesses), while for low proportions of object accesses, Marvin’s overhead was minimal (e.g., 10% overhead for 1% object accesses). PCMark’s 15% overhead indicates that the real app workloads represented by PCMark have low proportions of object accesses.

Although these overheads are already reasonable, they could be improved with optimizations. As noted in Section 7, deeper compiler integration would let Marvin reduce overhead by performing object access interposition less frequently.

Code size overhead of object access interposition. The ARM64 instructions added by Marvin’s object access interposition also increase the size of compiled native code. To measure the increase in code size, we compared the compiled Android framework libraries generated by Marvin to the framework libraries on unmodified Android. Marvin increased the total size of the ARM64 framework libraries (in the `/system/framework/arm64` and `/system/framework/oat/arm64` directories on the Android filesystem) from 117 MB to 292 MB. This code size overhead, while relatively high, could be reduced significantly with deeper compiler integration (Section 7).

CPU utilization overhead of heap walks. Our Marvin prototype performs the heap walks required for working set es-

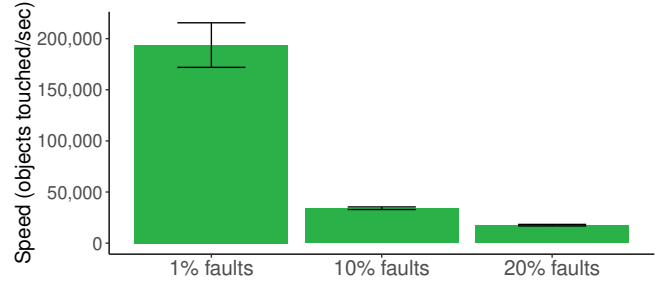


Figure 10: Speed of a benchmark app as it touches objects in its heap with different fractions of reclaimed objects.

timation by invoking the concurrent garbage collector and piggybacking off its heap walk. Our prototype performs a heap walk every 5 seconds when an app is in the foreground and every 30 seconds for an app in the background. In theory, this periodic invocation of the garbage collector across multiple MRT instances could add CPU utilization overhead. In practice, when running multiple apps in the background, we found that the difference between Marvin’s and unmodified Android’s CPU utilization was negligible, likely because GC invocations were so infrequent for background apps.

Overhead of faulting in objects. When an app first accesses a reclaimed object, Marvin must fault it in from disk, adding significant latency to that initial access. Marvin’s default policy eagerly restores all objects when an app moves to the foreground, trading off a longer transition delay for a guarantee that object-faulting latencies will never block the app’s UI thread once it is in the foreground. The added transition delay is proportional to the amount of reclaimed memory and the restoration rate. Anecdotally, our prototype restored memory at about 100 MB/s, but we believe that rate could improve to around 260 MB/s with optimization. (The slower rate is the disk read speed of the C++ standard library implementation used by ART, while the faster rate is the speed of the standard library implementation linked by the Android standalone toolchain.) Object faults may occur for background apps, but the Java working sets of apps in the background are generally quite small (less than 4MB for all commercial apps in our test set), so we expect object faulting to happen infrequently in practice.

We nonetheless studied the effect of object faulting on performance, to understand how Marvin would perform in situations where different policies or workloads result in more object faulting. Figure 10 shows the effect of object faulting on a heap-walking benchmark app as it touches different fractions of reclaimed objects. The app looped over a set of 4KB arrays, reading five member variables of each array, and measured the speed of traversing the objects. Each bar shows the mean and standard deviation of five measurements, and the device’s disk cache was cleared before each run. As expected,

there was an inverse relationship between heap-walking speed and fraction of faults; for instance, speed dropped by 49% as the fraction of faults increased from 10% to 20%.

9 Related Work

Several recent systems provide swapping for mobile platforms but focus on page-granularity rather than object-granularity swapping. SmartSwap [43] predicts which apps are unlikely to be used and swaps out pages from those apps ahead-of-time. A2S [22] takes the opposite approach; it avoids swapping out pages from unused apps, since their pages will be freed anyways when they are terminated. MARS [19] optimizes Linux swapping to improve performance on flash storage devices. It disables garbage collection in background apps and reclaims memory from those apps. DR. Swap [42] uses NVRAM rather than flash storage to store swapped-out pages and satisfies reads by reading directly from NVRAM. Choi et al. [6] improve the performance of an in-memory file system by co-designing the swap mechanism to minimize I/O.

The Linux kernel includes a daemon, *kswapd*, which frees unused pages in the background to maintain a reserve pool of unallocated memory [18]. Like Marvin, *kswapd* proactively checkpoints unused memory, but unlike Marvin, *kswapd* reclaims pages when it checkpoints them. As a result, *kswapd* is limited in how much memory it can checkpoint ahead-of-time—keeping a large proportion of memory checkpointed and reclaimed would make that memory unusable and shrink the device’s effective memory footprint.

Liang et al. [24] present FAST, an Android memory management system that modifies *kswapd* to improve its suitability for Android. FAST changes *kswapd* to prioritize reclaiming pages from apps in the background. It also identifies a mismatch between the large reclamation sizes of *kswapd* and the small sizes of typical Android allocations, and it includes a predictor to determine the reclamation size based on workload patterns. Like FAST, Marvin avoids reclaiming memory from the foreground app, but Marvin differs in its runtime-level memory management and its decoupling of checkpointing and reclamation.

A significant body of work examines the issue of providing persistent memory for object-oriented languages [1, 7, 26, 31, 33, 40]. These systems checkpoint objects to disk or non-volatile memory, but they do so to ensure safety in the face of failures rather than swapping out unused memory. As a result, they focus on supporting transactional programming models that provide strong guarantees under failure [7, 31] and on implementing crash-safe garbage collection [7, 40] rather than on maximizing the number of apps that can run concurrently.

SSDAlloc [2] is a persistent memory system that, like Marvin, is motivated by the goal of helping apps with large memory footprints avoid memory pressure. Unlike Marvin’s runtime-level object faulting and working set estimation, SSDAlloc allocates objects in separate virtual pages and uses

the existing virtual memory system to estimate the working set and trigger its object fault handler.

Like Marvin, the bookmarking collector [20] aims to improve the performance of Java apps in memory-constrained environments. It assumes that the OS uses a traditional page-level swapping mechanism and focuses on letting the garbage collector run without unnecessary swapping. It conservatively stores approximate reachability information (bookmarks) that is used during garbage collection, whereas Marvin stores exact reachability information (stubs). The BMX garbage collector [15] also uses stubs to avoid expensive object accesses, but in the context of a distributed persistent object store.

Other recent work on garbage collection focuses on co-designing the GC and runtime to manage software caches more efficiently [30], co-designing the GC and virtual memory manager to improve performance [41], measuring the effect of GC on scalability [16], and designing GCs or memory managers for domains such as big data systems [17, 29].

With multiple runtimes managing their own memory on top of a single operating system, Android’s architecture resembles that of a virtual machine manager, where multiple guest operating systems run on top of a hypervisor. Marvin’s runtime–OS cooperation is analogous to that between guest OS and hypervisor in the VMware ESX server [37], which uses a balloon driver to induce guest OSES to reclaim memory.

Wright et al. [39] present a system in which the architecture and Java runtime are co-designed for improved memory access performance. It features hardware modifications that allow an object-addressed CPU cache and an in-cache GC.

10 Conclusion

Users of mobile devices expect to use apps and switch between apps with low latency. As mobile apps have become more memory-hungry, device RAM capacities have not kept pace, and traditional swapping mechanisms cannot meet user latency expectations. Marvin overcomes this challenge with a novel runtime-level swapping mechanism that accurately estimates working sets, moves disk I/O off the allocation critical path, and avoids unnecessary swapping during garbage collection. As our experiments demonstrate, Marvin lets more apps run simultaneously and reclaims memory faster than unmodified Android while adding reasonable overhead. The source code for our Marvin prototype and experiments is available at <https://github.com/UWSysLab>.

Acknowledgments

We thank our anonymous reviewers and our shepherd, Malte Schwarzkopf, for their insightful feedback. We also thank Sandy Kaplan for editing and writing advice and Ali Razeen for help with Android troubleshooting. This work was funded in part by Futurewei.

References

- [1] Malcolm Atkinson and Ronald Morrison. Orthogonally persistent object systems. *The VLDB Journal*, 4(3):319–402, July 1995.
- [2] Anirudh Badam and Vivek S. Pai. Ssdalloc: Hybrid ssd/ram memory management made easy. In *Proceedings of the 8th USENIX Symposium on Networked Systems Design and Implementation (NSDI '11)*, 2011.
- [3] UL Benchmarks. Pemark for android. <https://benchmarks.ul.com/pemark-android>. Accessed: 2019-1-9.
- [4] Kofi Amankwah Boamah. iPhone on-board RAM, July 2017. https://www.researchgate.net/figure/Phone-on-board-RAM-From-figure-8-it-is-clear-that-Apple-either-maintains-the-iPhone_fig1_319307164.
- [5] Bumptech. Glide v4: Fast and efficient image loading for android. <https://bumptech.github.io/glide/>. Accessed: 2018-11-28.
- [6] J. Choi, J. Ahn, J. Kim, S. Ryu, and H. Han. In-memory file system with efficient swap support for mobile smart devices. *IEEE Transactions on Consumer Electronics*, 62(3):275–282, August 2016.
- [7] Joel Coburn, Adrian M. Caulfield, Ameen Akel, Laura M. Grupp, Rajesh K. Gupta, Ranjit Jhala, and Steven Swanson. Nv-heaps: Making persistent objects fast and safe with next-generation, non-volatile memories. In *Proceedings of the Sixteenth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS XVI*, pages 105–118, New York, NY, USA, 2011. ACM.
- [8] Peter J Denning and Stuart C Schwartz. Properties of the working-set model. *Communications of the ACM*, 15(3):191–198, 1972.
- [9] Android Documentation. Activity. <https://developer.android.com/reference/android/app/Activity>, 2018. Accessed: 2018-11-28.
- [10] Android Documentation. Caching bitmaps. <https://developer.android.com/topic/performance/graphics/cache-bitmap>, 2018. Accessed: 2018-11-30.
- [11] Android Documentation. Manage your app’s memory. <https://developer.android.com/topic/performance/memory>, 2018. Accessed: 2018-11-28.
- [12] Android Documentation. Saving ui states. <https://developer.android.com/topic/libraries/architecture/saving-states>, 2018. Accessed: 2018-11-28.
- [13] D. R. Engler, M. F. Kaashoek, and J. O’Toole, Jr. Exokernel: An operating system architecture for application-level resource management. In *Proceedings of the Fifteenth ACM Symposium on Operating Systems Principles, SOSP '95*, pages 251–266, 1995.
- [14] Umar Farooq and Zhijia Zhao. Runtimedroid: Restarting-free runtime change handling for android apps. In *Proceedings of the 16th Annual International Conference on Mobile Systems, Applications, and Services, MobiSys '18*, pages 110–122, 2018.
- [15] Paulo Ferreira and Marc Shapiro. Garbage Collection and DSM Consistency. In *Proceedings of the 1st USENIX Conference on Operating Systems Design and Implementation (OSDI '94)*, pages 229–241, Monterey CA, USA, United States, 1994.
- [16] Lokesh Gidra, Gaël Thomas, Julien Sopena, and Marc Shapiro. Assessing the scalability of garbage collectors on many cores. In *Proceedings of the 6th Workshop on Programming Languages and Operating Systems (PLOS '11)*, 2011.
- [17] Ionel Gog, Jana Giceva, Malte Schwarzkopf, Kapil Vaswani, Dimitrios Vytiniotis, Ganesan Ramalingam, Manuel Costa, Derek G. Murray, Steven Hand, and Michael Isard. Broom: Sweeping out garbage collection from big data systems. In *15th Workshop on Hot Topics in Operating Systems (HotOS XV)*, Kartause Ittingen, Switzerland, 2015.
- [18] Mel Gorman. An investigation into the theoretical foundations and implementation of the linux virtual memory manager, 2003.
- [19] Weichao Guo, Kang Chen, Huan Feng, Yongwei Wu, Rui Zhang, and Weimin Zheng. Mars: Mobile application relaunching speed-up through flash-aware page swapping. *IEEE Transactions on Computers*, 65(3):916–928, March 2016.
- [20] Matthew Hertz, Yi Feng, and Emery D. Berger. Garbage collection without paging. In *Proceedings of the 2005 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '05*, pages 143–153, 2005.
- [21] Tyler Kieft. Building a better instagram app for android. <https://instagram-engineering.com/building-a-better-instagram-app-for-android-c08f973662b>, 2014. Accessed: 2018-11-9.

- [22] Sang-Hoon Kim, Jinkyu Jeong, and Jin-Soo Kim. Application-aware swapping for mobile systems. *ACM Trans. Embed. Comput. Syst.*, 16(5s):182:1–182:19, September 2017.
- [23] Ian M. Leslie, Derek McAuley, Richard Black, Timothy Roscoe, Paul Barham, David Evers, Robin Fairbairns, and Eoin Hyden. The design and implementation of an operating system to support distributed multimedia applications. *IEEE Journal on Selected Areas in Communications*, 14(7):1280–1297, September 1996.
- [24] Yu Liang, Jinheng Li, Rachata Ausavarungnirun, Riwei Pan, Liang Shi, Tei-Wei Kuo, and Chun Jason Xue. Acclaim: Adaptive memory reclaim to improve user experience in android systems. In *Proceedings of the 2020 USENIX Annual Technical Conference*, 2020.
- [25] Michelle Meyers. Android inches ahead of windows as most popular os. CNET, April 2017. <https://www.cnet.com/news/android-most-popular-os-beats-windows-statcounter/>.
- [26] J. Eliot B. Moss. Design of the mneme persistent object store. *ACM Trans. Inf. Syst.*, 8(2):103–139, April 1990.
- [27] Mike Nakhimovich. Improving startup time in the nytimes android app. <https://open.blogs.nytimes.com/2016/02/11/improving-startup-time-in-the-nytimes-android-app/>, 2016. Accessed: 2018-11-9.
- [28] Randy Nelson. The size of iphone’s top apps has increased by 1,000% in four years. Sensor Tower, Jun 2017. <https://sensortower.com/blog/ios-app-size-growth>.
- [29] Khanh Nguyen, Lu Fang, Guoqing Xu, Brian Demsky, Shan Lu, Sanazsadat Alamian, and Onur Mutlu. Yak: A high-performance big-data-friendly garbage collector. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI ’16)*, 2016.
- [30] Diogenes Nunez, Samuel Z. Guyer, and Emery D. Berger. Prioritized garbage collection: Explicit gc support for software caches. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2016*, pages 695–710, New York, NY, USA, 2016. ACM.
- [31] James O’Toole, Scott Nettles, and David Gifford. Concurrent compacting garbage collection of a persistent heap. In *Proceedings of the Fourteenth ACM Symposium on Operating Systems Principles, SOSP ’93*, pages 161–174, New York, NY, USA, 1993. ACM.
- [32] Anshu Rustagi. How we improved our android app “cold start” time by 28%. <https://redfin.engineering/how-we-improved-our-android-app-cold-start-time-by-28-a722e231314a>, 2018. Accessed: 2018-11-9.
- [33] Vivek Singhal, Sheetal V. Kakkad, and Paul R. Wilson. Texas: An efficient, portable persistent store. In Antonio Albano and Ron Morrison, editors, *Persistent Object Systems*, pages 11–33, London, 1993. Springer London.
- [34] Facebook Open Source. Fresco. <https://frescolib.org/>. Accessed: 2018-11-28.
- [35] StackExchange. Creating and enabling an internal storage swap partition on rooted android kitkat. <https://android.stackexchange.com/a/89030>. Accessed: 2019-4-4.
- [36] StackOverflow. ios app maximum memory budget. <https://stackoverflow.com/a/15200855>. Accessed: 2019-1-9.
- [37] Carl A. Waldspurger. Memory resource management in vmware esx server. In *Proceedings of the 5th Symposium on Operating Systems Design and Implementation (OSDI ’02)*, 2002.
- [38] David A. Wheeler. SLOccount, 2013. <http://www.dwheeler.com/sloccount/>.
- [39] Greg Wright, Matthew L. Seidl, and Mario Wolczko. An object-aware memory architecture. Technical report, Sun Microsystems, Inc., Mountain View, CA, USA, 2005.
- [40] Mingyu Wu, Ziming Zhao, Haoyu Li, Heting Li, Haibo Chen, Binyu Zang, and Haibing Guan. Espresso: Brewing java for more non-volatility with non-volatile memory. In *Proceedings of the Twenty-Third International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS ’18*, pages 70–83, New York, NY, USA, 2018. ACM.
- [41] Ting Yang, Emery D. Berger, Scott F. Kaplan, and J. Eliot B. Moss. Cramm: virtual memory support for garbage-collected applications. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation (OSDI ’06)*, 2006.
- [42] K. Zhong, X. Zhu, T. Wang, D. Zhang, X. Luo, D. Liu, W. Liu, and E. H.-M. Sha. Dr. swap: Energy-efficient paging for smartphones. In *2014 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*, pages 81–86, Aug 2014.
- [43] Xiao Zhu, Duo Liu, Kan Zhong, Jinting Ren, and Tao Li. Smartswap: High-performance and user experience

friendly swapping in mobile systems. In *Proceedings of the 54th Annual Design Automation Conference 2017*,

DAC '17, pages 22:1–22:6, New York, NY, USA, 2017. ACM.