

Archipelago: Trading Address Space for Reliability and Security

Vitaliy B. Lvin, Gene Novark, Emery D. Berger, Benjamin G. Zorn[†]

Dept. of Computer Science
University of Massachusetts
Amherst, MA 01003

[†]Microsoft Research
One Microsoft Way
Redmond, WA 98052

ABSTRACT

Memory errors are a notorious source of security vulnerabilities that can lead to service interruptions, information leakage and unauthorized access. Because such errors are also difficult to debug, the absence of timely patches can leave users vulnerable to attack for long periods of time. A variety of approaches have been introduced to combat these errors, but these often incur large runtime overheads and generally abort on errors, threatening availability.

This paper presents Archipelago, a runtime system that takes advantage of available address space to substantially reduce the likelihood that a memory error will impact program execution. Archipelago randomly allocates heap objects far apart in virtual address space, effectively isolating each object from buffer overflows. Archipelago also protects against dangling pointer errors by preserving the contents of freed objects after they are freed. Archipelago thus trades *virtual* address space—a plentiful resource on 64-bit systems—for significantly improved program reliability and security, while limiting *physical* memory consumption by tracking the working set of an application and compacting cold objects. We show that Archipelago allows applications to continue to run correctly in the face of thousands of memory errors. Across a suite of server applications, Archipelago’s performance overhead is 6% on average (between -7% and 22%), making it especially suitable to protect servers that have known security vulnerabilities due to heap memory errors.

1. INTRODUCTION

Memory errors in C and C++ programs continue to be a significant problem. They are hard to debug and often easy to exploit. Memory-based attacks are an effective way to compromise Internet servers, either by crashing them, which causes service interruptions and data loss, or by making them execute arbitrary code. Because these bugs are difficult to debug, it can take weeks before even critical errors are repaired [31], leaving applications vulnerable to attack.

A variety of approaches have been developed to help programmers avoid memory errors. These approaches can be roughly classified into three categories: testing tools, garbage collectors, and compiler-based tools. Testing tools, such as Valgrind [22, 29] and Purify [14], impose performance overheads that only make their use feasible during testing. Conservative garbage collectors [6] only protect against dangling pointer errors and provide no protection against buffer overflows. Compiler-based approaches [1, 2, 9, 11, 17, 21, 24, 32, 34] typically incur unacceptably-large runtime overheads or require programmer intervention, and also require source code, which may not be available. They also generally abort program execution in response to memory errors, reducing availability and leaving systems vulnerable to denial-of-service attacks.

Contributions: This paper presents Archipelago, a runtime system that significantly improves the resilience of applications to heap-based memory errors.¹ Archipelago treats heap objects as individual islands, surrounded by stretches of unused address space. On modern architectures, especially 64-bit systems, virtual address space is a plentiful resource. Archipelago trades this plentiful resource for a high degree of *probabilistic memory safety* [3]; in other words, Archipelago can use available *virtual* memory to significantly increase the likelihood that a program will run correctly in the face of memory errors.

To control *physical* memory consumption, Archipelago leverages the following key insight: once the distance between objects crosses a certain threshold, each page will hold exactly one (small) object. At this point, additional address-space expansion is free: the virtual memory system does not need to allocate physical frames for unused address space between objects. Archipelago takes advantage of this insight and directly allocates one object per page, leaving the virtual address space between objects uncommitted. It further limits physical memory consumption by selectively compacting pages of the heap that are infrequently used.

The class of applications that are most sensitive to memory errors and associated security vulnerabilities are servers: they are attractive, high-value targets that are connected directly to the Internet. We show that Archipelago can provide high levels of safety and reliability for this class of applications. We show that Archipelago can let applications

¹An *archipelago* is an expanse of water with many scattered islands, such as the Aegean Sea.

run even in the face of thousands of memory errors, while keeping performance impact to acceptably-low levels. Archipelago slows down execution of a range of server applications by just 6% on average (from -7% to 22%). This modest performance impact makes Archipelago a realistic approach to protect deployed server applications against known and unknown heap-based security vulnerabilities.

While the primary focus of this work is using available memory to increase the availability and security of networked server applications, Archipelago can also *detect* memory errors. In this mode, Archipelago is both more efficient and more thorough than standard memory debugging tools. Because it generally allows programs to run despite memory errors, it can detect multiple heap overflows in a single execution, rather than halting on the first error.

The rest of the paper is organized as follows. Section 2.2 reviews operating system support for virtual memory, and explains probabilistic memory safety. Section 3 describes the software architecture of Archipelago in detail. Section 4 evaluates the effectiveness of Archipelago at withstanding memory errors and measures its overhead. Section 5 describes how Archipelago detects buffer overflows, and compares its use in this mode to two widely-used memory debuggers. Section 6 surveys related work, Section 7 discusses future directions, and Section 8 concludes.

2. BACKGROUND

2.1 Virtual Memory

Because Archipelago makes extensive use of operating system support for virtual memory management that may not be familiar, we define some key terms and concepts here.

A key distinction is that between virtual and physical memory. Virtual memory refers to the full addressable range of memory. This range does not necessarily correspond to the architecture’s word size—on x86-64 architectures, the addressable range is 48 bits (i.e., 2^{48} bytes). Operating systems map virtual memory to available physical memory. On 64-bit systems, virtual memory is plentiful while physical memory is in relatively short supply (e.g., on the order of 1–8 gigabytes (2^{30} – 2^{33}) bytes).

Virtual memory is divided into pages that are typically 4K chunks. Pages can be in three states: *unmapped*, *reserved*, and *committed*. An unmapped page is not available for use by the process, and access to it triggers a segmentation violation.

When a process obtains a page from the system via `mmap`, the virtual address range is *reserved* so that a subsequent call to `mmap` is guaranteed to return virtual memory from a different range. However, a reserved page does not initially have an associated physical page frame.

When a reserved page is touched for the first time, the page is *committed*: a physical page frame is allocated and associated with the virtual page. The kernel initializes all page contents to zero when they are first touched. Subsequent touches do not result in any page faults unless, due to memory pressure, the page is *evicted* to disk. In this case, the page’s contents are generally written to the disk, and then

the page is decommitted (but remains reserved). A subsequent touch triggers a page fault, and the kernel will fill the page with the contents previously saved on disk.

Operating systems allow programmers to control the state of pages via the `madvise` system call. A Unix application can invoke `madvise(MADV_FREE)` to inform the kernel that a range of pages is available to be reclaimed, and that there is no need to write the contents to disk. This call thus decommits a page’s physical frame, making it available for reuse by the system. Archipelago makes use of `madvise` to limit its physical memory footprint, as Section 3.1 describes. `madvise` can also be used to provide hints to guide the virtual memory manager’s page replacement algorithm, a feature that Archipelago also uses.

Additionally, an application can protect access to a page so that accesses trigger a signal, even if the page has been committed. For example, an application can invoke `mprotect(..., PROT_NONE)` on a range of pages: future attempts to read, write, or execute memory on this page will trigger a segmentation violation. By installing a custom signal handler to handle these segmentation violations, an application can selectively intercept reads or writes to particular pages. Archipelago uses these memory protection calls to let it perform compaction of cold objects (see Section 3.2).

2.2 Probabilistic Memory Safety

The motivation for our work comes from the ideas of *infinite heaps* and *probabilistic memory safety* originally introduced by Berger and Zorn [3].

An *infinite heap memory manager* is an ideal, unrealizable runtime system that allows programs containing memory errors to execute soundly and to completion. In such a system, the heap area is infinitely large and can never be exhausted. All objects are allocated fresh, infinitely far away from each other, and are never deallocated.

Because every object is infinitely far away from any other object, buffer overflows become benign, and dangling pointers also vanish since objects are never deallocated or reused. A portable correct C program cannot tell the difference between an infinite heap memory manager and a normal allocator, while a program containing memory errors would execute correctly for reasons outlined above, as long as it does not contain uninitialized reads.

Of course, it is impossible to build a true infinite heap memory manager. However, one can approximate its behavior by using an *M-heap*—a heap that is M times larger than needed. By placing objects uniformly randomly across an M -heap, we get the expected minimum separation between any two objects of $M - 1$ objects, and therefore overflows that are smaller become benign with high probability. By *randomizing* the choice of freed objects to reuse, we minimize the likelihood of recently freed objects being overwritten, and therefore of a malignant dangling pointer error. This heap thus provides *probabilistic memory safety*, a probabilistic guarantee that memory errors occurring in the program are benign during its execution.

```

1 void * malloc (size_t size) {
2     if (size <= PAGE_SIZE) {
3         //object fits on a page
4         //obtain random page from the pool
5         void *page = getRandomPage();
6     }
7     if (page == NULL) {
8         //object doesn't fit on the page
9         //or pool is full
10        //mmap memory directly
11        void *pages =
12            mmap(roundUpToPageSize(size),
13                MAP_ANONYMOUS);
14    }
15    if (page == NULL) {
16        //mmap failed
17        return NULL;
18    }
19    //add coloring
20    void *ptr =
21        getRandomColoring(page, size);
22    //register page(s) as part
23    //of working set
24    registerActivePages(page, ptr, size);
25    return ptr;
26 }

```

Figure 1: Pseudo-code for Archipelago’s malloc.

In an M -heap, the likelihood of no live objects being overwritten by an overflow N objects in size is $(1 - \frac{1}{M})^N$ [3].

Based on this formula, it is clear that one way to increase the probability of correct execution in the presence of memory error is to make the *heap expansion factor* (M) large. For example, $M = 100$ yields a 99% probability that a buffer overflow smaller or equal to the size of an object will be benign. It is impractical, however, to run DieHard system with large values of M because of its correspondingly large physical memory consumption (see Section 4).

Archipelago achieves these probabilistic guarantees against buffer overflows while consuming only a correspondingly large amount of *virtual* memory. It effectively controls physical memory consumption and provides lower CPU overheads than a comparably-sized DieHard heap, as Sections 4.2 and 4.3 show.

3. ARCHIPELAGO ARCHITECTURE

Archipelago consists of two parts: a randomizing memory allocator and a cold storage module, which controls the overall physical memory consumption of the program. These parts are compiled into a dynamically-linked library that, when pre-loaded before an executable, replaces standard memory management routines, such as `malloc` and `free`, with calls to the Archipelago allocator.

3.1 Object-Per-Page Allocator

Key to Archipelago’s protection from memory errors is its object-per-page memory allocator. It is constructed using

```

1 void free (void * ptr) {
2     //retrieve size
3     size_t size = getObjectSize(ptr);
4     //get first page
5     void *page = getStartPage(ptr);
6     //unregister pages being deleted
7     unregisterActivePages(page, ptr, size);
8     //discard pages
9     //that have been compacted
10    discardCompactedPages(page, ptr, size);
11    if (size <= PAGE_SIZE) {
12        //object fits on page:
13        //discard contents
14        madvise(page, MADV_FREE);
15    } else {
16        //object doesn't fit on page:
17        //unmap it
18        munmap(page,
19            roundUpToPageSize(size));
20    }
21 }

```

Figure 2: Pseudo-code for Archipelago’s free.

the Heap Layers infrastructure [5]. As implied by its name, the object-per-page allocator places each allocated object on a separate virtual memory page. It reserves (but does not commit) a large fraction of the address space using `mmap`, and uses it as a pool from which to draw pages to satisfy allocation requests. Figures 3.1 and 3.1 present pseudo-code for `malloc` and `free`.

The size of the pool of available pages is a parameter to Archipelago (defaulting to 512 megabytes) that represents the trade-off between the protection Archipelago provides and its virtual memory consumption. A larger pool will provide more robust protection against errors, but at the cost of increased virtual memory consumption. Note that in case of memory pressure, the virtual memory manager will reclaim all committed but unused pages in the pool first, making the footprint of the application independent of the pool size.

Allocation: Objects are placed on pages randomly chosen from the pool (Figure 3.1, line 5). The object-per-page allocator uses a bit array to distinguish between used (allocated) and unused pages in the pool, and probes in the bit array to perform this random selection. In order to bound the expected number of probes to find an empty page, the object-per-page allocator always keeps the pool no more than half full. This strategy bounds the worst-case expected number of probes to a small constant (2).

Notice that since pages in the pool are allocated randomly, no locality of reference exists between pages in the pool. We give a hint to the virtual memory manager that no locality exists and that it should not prefetch pages within the pool using `madvise` (not shown in the code). Together with `mmap`, `madvise` ensures that pages are not instantiated in the physical memory until they are actually needed.

To reduce cache conflicts, Archipelago uses *colors* to place objects on pages: objects are placed at random offsets on pages, taking care to keep objects within their pages' boundaries (lines 20–21). Coloring helps reduce L2 misses due to cache conflicts and thus improves performance (we do not report these results here due to space limitations).

Deallocation: When an object smaller than a page in size is deleted, the object-per-page allocator marks the page as free (Figure 3.1, lines 5–10). Moreover, it instructs the virtual memory manager using `madvise` to discard the contents of the page without writing them to disk, therefore reducing the overall runtime overhead of the system due to page eviction (line 14).

Large objects: Objects that do not fit on a single page are treated specially by the object-per-page allocator. Archipelago currently does not search for ranges of free pages in the pool but instead allocates memory directly using `mmap` (Figure 3.1, lines 7–13). When the memory pool gets more than half full, all objects are allocated via `mmap` to avoid costly probes for free pages in the pool.

Because current Linux kernels randomize locations of memory-mapped objects in the address space, the object-per-page allocator need not take further action. When an object that was allocated using `mmap` is freed, its memory is immediately released back to the operating system using `munmap` (Figure 3.1, lines 18–19).

3.2 Exploiting Working Sets

Running programs with the object-per-page allocator alone would consume so much physical memory that it would be impractical for deployed programs. In order to limit its physical memory consumption, Archipelago relies on the observed temporal locality of memory accesses in most programs, or the so-called *working set hypothesis*. A program at any given time has a *working set*, a (hopefully) small subset of all live objects on which the program is actively operating.

The notion of a working set is extensively used in the virtual memory managers [8], which attempt to keep just the working sets of running programs in memory while storing rarely used data in secondary storage.

Archipelago follows a similar design: it tracks the working set of a program, moves objects not in the working set into a more compact representation, and returns the physical pages they occupied to the OS.

Operating systems typically rely on hardware-managed dirty and reference bits that give them precise information about which pages are being used. While a similar approach can be implemented in user-space with memory protection mechanisms, the cost of such an approach would be prohibitively high compared to the cost of making a mistake and compacting a page that is in the working set.

Instead, Archipelago uses a cheap approximation of the working set. Archipelago keeps all the pages occupied by live objects in a bounded FIFO queue. In our current implementation, the size of the FIFO queue is fixed at startup

```

1 void deflate (void *page) {
2     // allocate space in cold store
3     void *coldStore = coldHeap.malloc(
4         hotPages[page]->getDataSize());
5     // move the data
6     memcpy(coldStore,
7         hotPages[page]->getDataStart(),
8         size);
9     // return physical page to OS
10    madvise(page, MADV_FREE);
11    // set trap on future accesses
12    mprotect(page, PROT_NONE);
13    // mark page as cold
14    coldPaged[page] = hotPages[page];
15    hotPages.remove(page);
16    // remember the location of the data
17    coldPages[page].
18        setColdStore(coldStore);
19 }
20
21 bool inflate (void *page) {
22     // check page was deflated before
23     if (!coldPages.containsKey(page))
24         return false;
25     // enable access to the page
26     mprotect(page, PROT_READ | PROT_WRITE);
27     // restore data
28     memcpy(coldPages[page].getStart(),
29         coldPages[page].getColdStore(),
30         coldPages[page].getSize());
31     // free the cold space
32     coldHeap.free(
33         coldPages[page].getColdStore());
34     // mark page as hot
35     hotPages[page] = coldPages[page]
36     coldPages.remove(page);
37     return true;
38 }
39
40 void sigsegv_handler(void *addr) {
41     if (!inflate(getPageStart(addr))) {
42         fprintf(stderr, "Overflow!\n")
43         exit(-1);
44     }
45 }

```

Figure 3: Pseudo-code for Archipelago’s compaction and uncompression routines.

time, either read in from an environment variable or defaulting to 5000 objects. Pages are added to the back of the queue at allocation time. As the queue becomes full, pages at the front of the queue are removed and compacted. Upon access to a compacted page, the page is restored and added to the end of the queue as well.

3.3 Cold Storage

Archipelago contains functionality that compacts pages not in the current working set, thus reducing its physical memory requirements. It uses an in-memory compaction mechanism that stores compacted objects in a separate heap man-

aged by the general-purpose Lea allocator [19].

When a page is compacted, its non-zero contents are copied out into this internal heap. Then, any accesses to the page are disabled by `mprotect`, so that Archipelago receives a protection violation signal the next time the application tries to access the page. Finally, the virtual memory manager is instructed using `madvise` not to write the page contents to disk.

Archipelago installs a custom signal handler to receive protection violation signals and restore objects back from cold storage. When the handler receives a signal, it first checks whether the application was trying to access a page in cold storage. If it was, the handler has to restore that page in order for the application to continue. The handler unprotects the page and restores the data on it from cold storage. It also places the page back on the queue of active pages, and frees the space used to hold the page’s data in cold storage. Control then passes back to the application, which can now safely continue.

While compacting pages imposes additional runtime overhead, it effectively controls physical memory overhead, as Section 4.3 shows.

4. EVALUATION

In our evaluation, we answer the following questions:

1. What is the runtime overhead of using Archipelago with server applications?
2. What is the memory overhead of using Archipelago with server applications?
3. How effective is Archipelago against both injected faults and real errors?

4.1 Experimental Methodology

We perform our evaluation on a quiescent dual-processor with 8 gigabytes of RAM. Each processor is a 4-core 64-bit Intel Xeon running at 2.33 Ghz and equipped with a 4MB L2 cache.

We compare Archipelago to the GNU C library, which uses a variant of the Lea allocator [19], and to DieHard, version 1.1. This version, available from the project website, is an adaptive variant that dynamically grows its heap [4], and so is more space-efficient than the original, published description [3].

One important caveat is that we run all experiments on a particular version of a recent Linux kernel, version 2.6.21-mm2. This kernel version uses a more sophisticated algorithm for managing physical memory pages that were initially used by applications, but then returned to the kernel. This *page laundering* process updates a number of kernel data structures and potentially writes the page’s contents to secondary storage. Linux kernel versions up to and including 2.6.21 launder pages *eagerly* whenever an application calls `madvise`. However, Linux version 2.6.21-mm2 launders pages *lazily*, waiting until more physical memory pages are

actually needed. In the absence of memory pressure, this policy improves our system’s performance on an allocation-intensive microbenchmark by a factor of two: `madvise` is on Archipelago’s normal deallocation path. Because of its performance advantages for ordinary workloads, we expect that this patch, or one similar to it, will be adopted in future versions of the Linux kernel.

4.2 Performance Overhead

To quantify the performance overhead of using Archipelago, we measure the performance of a range of server applications running with and without Archipelago. All observed variances were below 1%. In our experiments, Archipelago uses a memory pool 512MB in size. We also compare performance against DieHard with two different heap multiplier values: 2 and 1024. The first multiplier provides performance and protection similar to the results reported in the original DieHard paper, while the second multiplier more closely approaches the level of protection that Archipelago achieves.

We use three different server applications: the *thttpd* web server, the *bftpd* ftp server, and an *openssh* server. For the first two, we record total throughput achieved with 50 simultaneous clients issuing 100 requests each. For the *openssh* server, we record the time it takes to perform authentication, spawn a shell, and disconnect, averaged over 10 runs.

We focus on the CPU impact of our benchmarks by performing all our experiments over the loop-back interface, so that any performance impact is not swamped by network latency. These measured runtime overheads are thus conservative estimates of the performance overhead one would see in practice.

Figure 4 presents the results of these experiments, normalized to GNU libc. These results show that Archipelago can protect servers without unduly sacrificing server performance. The performance overhead we observe is generally less than 20%. *thttpd* running with Archipelago repeatably performs better than with GNU libc; we do not yet fully understand why.

To evaluate the worst-case overhead one could expect for Archipelago, we also measure the performance impact of Archipelago on a well-known, extremely allocation-intensive benchmark, *espresso*. *espresso* allocates and deallocates approximately 1.5 million objects in less than a second. This allocation rate far exceeds that of a typical server application. In our experiments, we run *espresso* with all four allocators we use in our server experiments. Compared to GNU libc, *espresso* runs 3.34, 7.24 and 7.32 times slower with DieHard-2, DieHard-1024 and Archipelago, respectively.

4.3 Space Overhead

We evaluate the additional memory consumption incurred by using Archipelago, and compare this to DieHard and GNU libc.

Figure 5 shows virtual memory consumption of three different servers in our experiments. Due to the fact that Archipelago preallocates a large memory pool at start-up, its virtual memory consumption is always high compared to GNU libc.

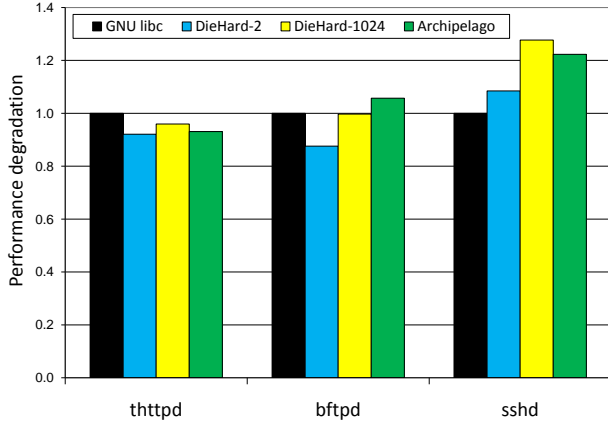


Figure 4: Performance across a range of server applications (smaller is better), normalized to GNU libc.

A large fraction—more than 70%—of that allocated space is never actually committed to memory. The high memory consumption of *bftpd* is explained by the fact that it forks three processes for every connected client. In our experiments, we use 50 simultaneous clients, and measure the total memory used by all *bftpd* processes, which includes the overheads of Archipelago in each *bftpd* process.

Figures 6(a) and 6(b) show resident memory consumption of *thttpd*, *bftpd*, and *sshd* during our experiments without and with memory pressure, respectively. We simulate memory pressure by locking pages in memory so that only about 512MB is usable by the entire system. Our experiments show that Archipelago uses less memory than DieHard-1024 and uses on average 3 to 5 times as much memory as GNU libc. This number is much lower for *thttpd* and *sshd*, which do not spawn multiple processes. It is important to note that memory consumption with Archipelago in the absence of memory pressure is artificially inflated, because Linux reclaims available pages only under memory pressure.

4.4 Avoiding Injected Faults

We evaluate the effectiveness of Archipelago in tackling memory errors by using two different types of fault injectors: an overflow injector and a dangling pointer injector. We inject faults into *espresso* running with GNU libc, DieHard and Archipelago. We perform all our injection experiments 100 times, and record the number of times that *espresso* produces correct output. Table 1 summarizes these results.

Buffer overflows: We perform three sets of experiments with the overflow injector. We inject 8-byte overflows with 0.01 probability, 4K overflows with 0.001 probability, and 8K overflows with 0.0001 probability. These probabilities correspond to thousands, hundreds, and tens of injected faults, respectively.

In this set of experiments, GNU libc crashes every time, as expected. Archipelago substantially outperforms both variants of DieHard across the range of overflow sizes and

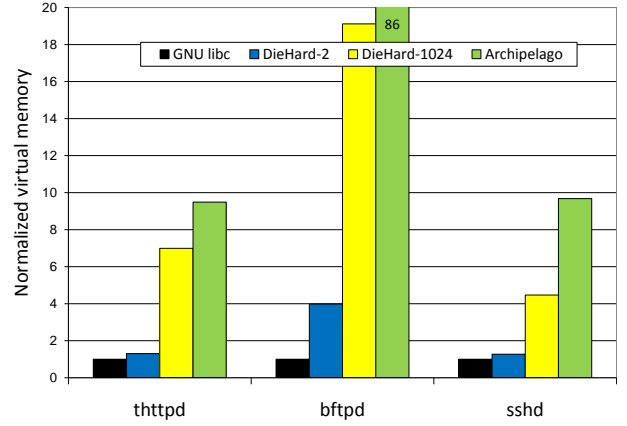


Figure 5: Virtual memory usage of GNU libc, DieHard, and Archipelago, across a range of server applications.

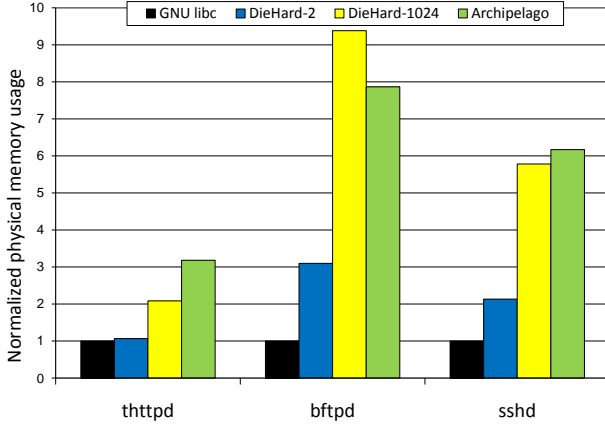
frequencies. With small and frequent overflows, Archipelago runs correctly every time. DieHard-1024 does reasonably well, running correctly 77% of the time, while DieHard-2 only runs correctly 29% of the time.

With large but infrequent overflows, Archipelago runs correctly 68% of the time. In this case, DieHard-1024 runs correctly only 23% of the time, while DieHard-2 crashes every time. Even in the worst case of large and reasonably frequent overflows, Archipelago lets *espresso* run correctly 42% of the time, while it only runs 2% of the time with DieHard-1024 (DieHard-2 crashes every time in this case).

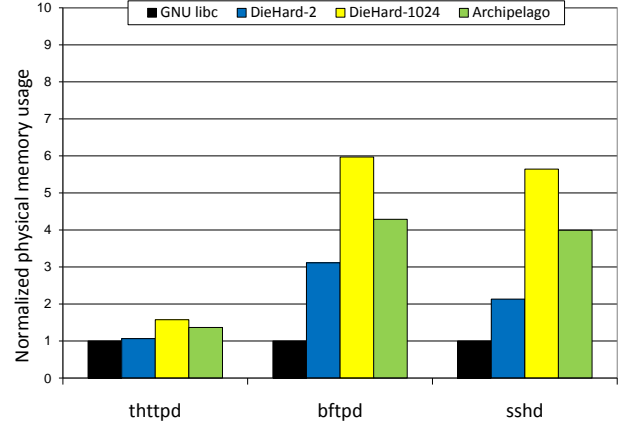
These results show that Archipelago provides excellent protection against buffer overflows and offers dramatic improvement over DieHard, even with an expansion factor of 1024.

Dangling pointers: Archipelago’s design goal was to limit the impact of buffer overflows, but it also provides a measure of protection against dangling pointers. To measure the impact of dangling pointers on runtime systems, we injected dangling pointer faults that free objects 5, 10 and 20 allocations early with probabilities 0.01, 0.001 and 0.0001, respectively.

These experiments show that, as expected, DieHard-1024 offers better protection from dangling pointer errors than Archipelago. This result is explained by the fact that DieHard-1024 has vastly more available object slots for reuse than Archipelago does. Archipelago has fewer potential slots to place new objects, since it only allows one object per page. Archipelago also instructs the operating system that all freed objects are available for the operating system to reuse at its discretion. If the operating system reuses a page, the original contents will be lost, and access through a dangling pointer to this data will trigger a fault. Nonetheless, Archipelago provides substantial protection against these errors, running correctly 29% of the time in the first experiment, 67% of the time in the second, and 98% in the third.



(a) Resident memory usage, without memory pressure.



(b) Resident memory usage, with memory pressure.

Figure 6: Resident memory usage with and without memory pressure. Under memory pressure, Linux quickly reclaims Archipelago’s uncommitted pages, making its physical memory consumption strictly lower than with DieHard-1024.

Injection experiments (% correct executions)				
espresso	GNU libc	DieHard-2	DieHard-1024	Archipelago
<i>buffer overflows</i>				
8 bytes, $p = 0.01$	0%	29%	77%	100%
8K, $p = 0.0001$	0%	0%	23%	68%
4K, $p = 0.001$	0%	0%	2%	42%
<i>dangling pointers</i>				
5 mallocs, $p = 0.01$	0%	8%	91%	29%
10 mallocs, $p = 0.001$	0%	75%	100%	67%
20 mallocs, $p = 0.0001$	0%	96%	100%	98%

Table 1: The performance of various runtime systems in response to injected memory errors (Section 4.4). Archipelago provides the best protection against overflows of all sizes and frequencies, and reasonable protection against dangling pointer errors (all executions fail with GNU libc).

4.5 Avoiding Real Buffer Overflows

To evaluate the effectiveness of Archipelago against real-life buffer overflows, we reproduce two well-known buffer overflow-based exploits: one in the *pine* mail reader, and the other in the *Squid* web cache proxy.

We reproduce an exploit in *pine* version 4.44. The exploit is a buffer overflow that can be triggered by a maliciously formed email message and causes *pine* to crash and fail to restart until the message is manually removed. When we place a maliciously formed message in a user’s mailbox, *pine* with GNU libc crashes whenever the user attempts to open a mailbox. However, when running with Archipelago, *pine* successfully opens the mailbox and performs all standard operations with messages in it, including the malicious message, without any user-noticeable slowdown.

We also test Archipelago’s ability to withstand a heap buffer overflow for the *squid* web cache. For version 2.3.STABLE5, a maliciously formed request causes a buffer overflow that corrupts heap meta-data (this causes GNU libc to terminate). When running with Archipelago, *squid* consistently handles the malicious request correctly, without crashing.

5. OVERFLOW DETECTION

Archipelago not only lets programs to run in the face of memory errors, but can also be used to detect and report these errors. Archipelago detects these errors using three separate approaches. First, Archipelago clears all memory pool pages before their first use. Because every page is initialized to zero, any non-zero value past the end of an object’s allocated space indicates an overflow. Archipelago scans the contents of a page past an allocated object on every **free**.

Second, Archipelago piggybacks buffer overflow detection onto page compaction, letting it discover overflows before an object is deallocated. Whenever a page is compacted, Archipelago recomputes the actual size of the object on that page by scanning it backwards until the first non-zero word. As above, any non-zero past the end indicates an overflow.

Finally, if an overflow touches a protected or unmapped page, Archipelago reports this as a heap overflow.

5.1 Evaluation

To evaluate Archipelago’s buffer overflow detection, we compare it to Valgrind (using the Memcheck tool [29]) and Electric Fence [25]. These tools differ substantially from Archipelago in their approaches. Valgrind uses heavyweight dynamic binary instrumentation to insert run-time checks; the Memcheck tool detects a wide range of memory errors (not just heap overflows) throughout program execution. Electric Fence is a debugging allocator that, like Archipelago, allocates heap objects on separate pages. It allocates three pages for every object: one page for the object itself (placed at the end), and a memory-protected page before and after the object. Electric Fence aborts whenever an overflow causes a memory protection fault.

We measure Archipelago’s overhead in overflow detection mode by running *espresso*, which effectively measures its worst-case performance. Archipelago is more efficient than Electric Fence (10.5x faster) and Valgrind (4.65x faster).

We then ran all three tools on *pine*, attempting to detect the buffer overflow we exploit in Section 4.5. Both Electric Fence and Valgrind successfully detect a single buffer overflow, but then abort the computation. However, Archipelago allows *pine* to safely continue execution and detects a second overflow. This experiment demonstrates Archipelago’s advantage over other tools: it can detect *multiple* heap overflows in a single run.

6. RELATED WORK

This section first discusses past work that exploits large address spaces, and then describes related work in the spheres of memory management, fault tolerance, and software engineering that address the problem of memory errors in C/C++ programs.

The advent of 64-bit processors sparked research in operating systems designed for large address spaces [7]. Druschel and Peterson point out that this address space is sufficiently large that it can be used to provide high performance protection and security by hiding processes from each other [13]. Anonymous RPC (ARPC) uses random placement of processes in a large address space to eliminate expensive hardware context switches on cross-domain RPC calls [33]. We are also leveraging a large address space, but instead of using the space to protect independent processes from each other, we are isolating individual objects from memory errors within the same process.

Archipelago builds on the ideas of Berger and Zorn’s DieHard system [3]. Like DieHard, Archipelago uses a randomized memory manager to provide protection from buffer overflows and dangling pointer errors. Unlike DieHard, Archipelago achieves high reliability by dramatically increasing the size of the address space and does not use replication. By exploiting both standard OS mechanisms and common program behavior, Archipelago provides greater resilience to buffer overflow errors with moderate and acceptable CPU and memory overhead.

Exterminator is another runtime system that, like DieHard, is based on randomized, overprovisioned heaps [23]. The focus of Exterminator is on automatic error detection and

correction based on accumulating data from multiple executions. While Archipelago can also be used for overflow detection, it is closer in spirit to DieHard, and unlike Exterminator, provides greater error tolerance without the requirement that errors first be detected.

A number of compiler-assisted approaches have been introduced to combat memory errors. Semantics provided by Archipelago to programs containing buffer overflows are similar to those of Rinard et al.’s Boundless Memory Blocks [27]. Because Boundless Memory Blocks uses a fixed-size LRU cache to store the values of out-of-bounds writes, accesses to out-of-bounds addresses are undefined if the object has been evicted from the cache. A number of other unsound approaches have been proposed [12, 28]. Dhurjati et al. use pool allocation to provide an efficient form of memory safety that guarantees that structure fields are referenced with the correct type. While they guarantee type-safety, there is no guarantee that the object the programmer had intended to access is correctly accessed [12]. Unlike this previous work, Archipelago provides a strong, quantifiable probabilistic guarantee that the intended program behavior will be preserved.

More traditional safe-C compilers [30, 21, 24] use modified versions of C and some combination of static analysis and dynamic checks to provide protection from memory errors. Cyclone [17, 30] augments C with an advanced type system to provide safe explicit memory management. CCured [21] inserts dynamic checks that ensure safety into the compiled program and uses static analysis to eliminate checks from places where memory errors cannot occur. CRED [24] only targets string buffer overflows, and inserts dynamic checks on memory accesses that use out-of-bounds pointers. All of these techniques are aimed at detecting memory errors and terminating the program in response. Archipelago, on the other hand, is aimed at avoiding memory errors and allowing the program to continue running correctly.

Like Archipelago, Rx can help avoid memory errors [26]. It performs periodic checkpointing of program execution, and when an error occurs, it re-runs the program from a checkpoint in a modified environment. In response to crashes, Rx pads allocations to avoid buffer overflows, and delays reuse of freed memory to prevent dangling pointers. Two fundamental limitations of Rx are that it only works with applications that allow replay, and cannot cope with errors that do not result in crashes. Archipelago does not suffer from either of these limitations.

Dangling pointer errors have been addressed in several ways in previous work. Dhurjati et al. employ a clever use of virtual memory page mapping and protection to allow them to detect dangling pointers at low cost [10]. While Archipelago also uses virtual memory protection, our focus is on providing resilience to buffer overflows with less emphasis on dangling pointers. Garbage collection is an alternative runtime system that provides safety from dangling pointer errors. The most commonly used garbage collector for C programs is Boehm-Demers-Weiser conservative garbage collector [6]. Unlike Archipelago, garbage collection provides no protection against buffer overflow errors. Garbage collection also imposes significant space and time overheads to achieve rea-

sonable performance [15].

Finally, a number of testing tools and debugging allocators [18, 20, 25] can aid programmers in debugging memory errors. Valgrind [22, 29] and Purify [14] use binary instrumentation or emulation to detect memory errors at run-time. They typically incur prohibitively high overhead both in terms of performance (up to 25X) and space (10X), making them only suitable during testing. Electric Fence [25] uses page protection to detect buffer overflow and dangling pointer errors. As we show in section 5, Electric Fence incurs high performance overhead and memory overhead, especially when used to detect dangling pointer errors. Unlike both of these systems, Archipelago can detect multiple heap buffer overflows in a single execution.

7. FUTURE WORK

There are a number of ways that the current Archipelago implementation can be improved. We intend to explore adaptively sizing the memory pool size to achieve the optimal trade-off between performance overhead and resilience to errors. Our current implementation has a static FIFO size, and we intend to investigate techniques to grow and shrink the FIFO size just as an OS virtual memory manager adapts working set size. Because our approach uses very large virtual address spaces with sparsely mapped pages, we will investigate how OS support for sparse page tables can improve Archipelago performance. Hardware TLBs have remained relatively small despite enormous growth in physical memory sizes over the last two decades. We anticipate that TLB designs that better accommodate large sparse virtual memories, such as those proposed by Huck and Hays [16], will significantly benefit Archipelago's performance.

8. CONCLUSION

Archipelago is a runtime system that provides protection from memory errors for unmodified C programs. It provides probabilistic protection from both buffer overflows and dangling pointer errors with high probability. Archipelago spreads objects far apart in the address space and randomizes the choice of freed objects to reuse, giving applications an illusion of infinite-size heap and protecting them from memory errors. It leverages the virtual memory subsystem of the underlying OS to efficiently provide a high level of memory safety to target programs at low cost.

We show analytically and empirically that Archipelago increases the resilience of programs to memory errors. Our evaluation shows that Archipelago is effective against real and injected memory errors. We show that it allows programs to correctly execute through hundreds and even thousands of memory errors, which is a significant improvement over current state-of-the-art systems. We also show how our system can be used to debug buffer overflow errors.

We further demonstrate that the overhead of using Archipelago is more than acceptable across a range of different server applications, both in terms of CPU performance and memory usage. We believe Archipelago is especially suitable for deployment in order to protect servers that have known security vulnerabilities due to heap memory errors.

9. REFERENCES

- [1] T. M. Austin, S. E. Breach, and G. S. Sohi. Efficient detection of all pointer and array access errors. In *Proceedings of the ACM SIGPLAN 1994 Conference on Programming Language Design and Implementation*, pages 290–301, New York, NY, USA, 1994. ACM Press.
- [2] D. Avots, M. Dalton, V. B. Livshits, and M. S. Lam. Improving software security with a C pointer analysis. In *ICSE '05: Proceedings of the 27th international conference on Software engineering*, pages 332–341, New York, NY, USA, 2005. ACM Press.
- [3] E. D. Berger and B. G. Zorn. DieHard: Probabilistic memory safety for unsafe languages. In *Proceedings of the 2006 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2006)*, pages 158–168, New York, NY, USA, 2006. ACM Press.
- [4] E. D. Berger and B. G. Zorn. Efficient probabilistic memory safety. Technical Report UMCS TR-2007-17, Department of Computer Science, University of Massachusetts Amherst, Mar. 2007.
- [5] E. D. Berger, B. G. Zorn, and K. S. McKinley. Composing high-performance memory allocators. In *Proceedings of the 2001 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2001)*, Snowbird, Utah, June 2001.
- [6] H.-J. Boehm and M. Weiser. Garbage collection in an uncooperative environment. *Software Practice and Experience*, 18(9):807–820, 1988.
- [7] J. S. Chase, H. M. Levy, M. J. Feeley, and E. D. Lazowska. Sharing and protection in a single-address-space operating system. *ACM Transactions on Computer Systems*, 12(4):271–307, Nov. 1994.
- [8] P. J. Denning. The working set model for program behaviour. *Communications of the ACM*, 11:323–333, 1968.
- [9] D. Dhurjati and V. Adve. Backwards-Compatible Array Bounds Checking for C with Very Low Overhead. In *Proceedings of the 2006 International Conference on Software Engineering (ICSE'06)*, Shanghai, China, May 2006.
- [10] D. Dhurjati and V. Adve. Efficiently detecting all dangling pointer uses in production servers. In *DSN '06: Proceedings of the International Conference on Dependable Systems and Networks (DSN'06)*, pages 269–280, Washington, DC, USA, 2006. IEEE Computer Society.
- [11] D. Dhurjati, S. Kowshik, and V. Adve. Safecode: enforcing alias analysis for weakly typed languages. In *Proceedings of the 2006 ACM SIGPLAN conference on Programming language design and implementation*, pages 144–157, New York, NY, USA, 2006. ACM Press.
- [12] D. Dhurjati, S. Kowshik, V. Adve, and C. Lattner. Memory safety without runtime checks or garbage collection. In *ACM SIGPLAN 2003 Conference on Languages, Compilers, and Tools for Embedded Systems (LCTES'2003)*, San Diego, CA, June 2003. ACM Press.
- [13] P. Druschel and L. L. Peterson. High-performance

- cross-domain data transfer. Technical Report TR 92-11, Dept. Comp. of Sc., U. of Arizona, Tucson, AZ (USA), Mar. 1992.
- [14] R. Hastings and B. Joyce. Purify: Fast detection of memory leaks and access errors. In *Proc. of the Winter 1992 USENIX Conference*, pages 125–138, San Francisco, California, 1991.
 - [15] M. Hertz and E. D. Berger. Quantifying the performance of garbage collection vs. explicit memory management. In *Proceedings of the 20th annual ACM SIGPLAN Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA)*, San Diego, CA, Oct. 2005.
 - [16] J. Huck and J. Hays. Architectural support for translation table management in large address space machines. In *ISCA '93: Proceedings of the 20th annual international symposium on Computer architecture*, pages 39–50, New York, NY, USA, 1993. ACM Press.
 - [17] T. Jim, J. G. Morrisett, D. Grossman, M. W. Hicks, J. Cheney, and Y. Wang. Cyclone: A safe dialect of C. In *Proceedings of the General Track: 2002 USENIX Annual Technical Conference*, pages 275–288, Berkeley, CA, USA, 2002. USENIX Association.
 - [18] M. Kharbutli, X. Jiang, Y. Solihin, G. Venkataramani, and M. Prvulovic. Comprehensively and efficiently protecting the heap. In *ASPLOS-XII: Proceedings of the 12th International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 207–218, New York, NY, USA, 2006. ACM Press.
 - [19] D. Lea. A memory allocator. <http://gee.cs.oswego.edu/dl/html/malloc.html>, 1997.
 - [20] Microsoft Corporation. Pageheap. <http://support.microsoft.com/kb/286470>.
 - [21] G. C. Necula, S. McPeak, and W. Weimer. CCured: type-safe retrofitting of legacy code. In *POPL '02: Proceedings of the 29th ACM SIGPLAN-SIGACT symposium on Principles of Programming Languages*, pages 128–139, New York, NY, USA, 2002. ACM Press.
 - [22] N. Nethercote and J. Fitzhardinge. Bounds-checking entire programs without recompiling. In *SPACE 2004*, Venice, Italy, Jan. 2004.
 - [23] G. Novark, E. D. Berger, and B. G. Zorn. Exterminator: automatically correcting memory errors with high probability. In *PLDI '07: Proceedings of the 2007 ACM SIGPLAN conference on Programming language design and implementation*, pages 1–11, New York, NY, USA, 2007. ACM Press.
 - [24] O. Ruwase and Monica S. Lam. A practical dynamic buffer overflow detector. In *Proceedings of the 11th Annual Network and Distributed System Security Symposium*, pages 159–169, Feb. 2004.
 - [25] B. Perens. Electric Fence v2.1. <http://perens.com/FreeSoftware/ElectricFence/>.
 - [26] F. Qin, J. Tucek, J. Sundaresan, and Y. Zhou. Rx: Treating bugs as allergies: A safe method to survive software failures. In *Proceedings of the Twentieth Symposium on Operating Systems Principles*, volume XX of *Operating Systems Review*, Brighton, UK, Oct. 2005. ACM.
 - [27] M. Rinard, C. Cadar, D. Dumitran, D. M. Roy, and T. Leu. A dynamic technique for eliminating buffer overflow vulnerabilities (and other memory errors). In *Proceedings of the 2004 Annual Computer Security Applications Conference*, Dec. 2004.
 - [28] M. Rinard, C. Cadar, D. Dumitran, D. M. Roy, T. Leu, and J. William S. Beebe. Enhancing server availability and security through failure-oblivious computing. In *Sixth Symposium on Operating Systems Design and Implementation*, San Francisco, CA, Dec. 2004. USENIX.
 - [29] J. Seward and N. Nethercote. Using Valgrind to detect undefined value errors with bit-precision. In *Proceedings of the USENIX'05 Annual Technical Conference*, Anaheim, California, USA, Apr. 2005.
 - [30] N. Swamy, M. Hicks, G. Morrisett, D. Grossman, and T. Jim. Experience with safe manual memory management in cyclone. *Science of Computer Programming*, 2006. Special issue on memory management. Expands ISMM conference paper of the same name. To appear.
 - [31] Symantec. Internet security threat report. <http://www.symantec.com/enterprise/threatreport/index.jsp>, Sept. 2006.
 - [32] W. Xu, D. C. DuVarney, and R. Sekar. An efficient and backwards-compatible transformation to ensure memory safety of C programs. In *SIGSOFT '04/FSE-12: Proceedings of the 12th ACM SIGSOFT twelfth international symposium on Foundations of software engineering*, pages 117–126, New York, NY, USA, 2004. ACM Press.
 - [33] C. Yarvin, R. Bukowski, and T. Anderson. Anonymous RPC: Low-latency protection in a 64-bit address space. In *Proceedings of the 1993 Summer USENIX Conference*, pages 175–186, 1993.
 - [34] S. H. Yong and S. Horwitz. Protecting C programs from attacks via invalid pointer dereferences. In *ESEC/FSE-11: 11th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 307–316, New York, NY, USA, 2003. ACM Press.