# Cichlid: Explicit physical memory management for large machines

*Simon Gerber, Gerd Zellweger, Reto Achermann,*
*Moritz Hoffmann, Kornilios Kourtis, Timothy Roscoe, Dejan Milojicic*[†]
*Systems Group, Department of Computer Science, ETH Zurich*      [†]*Hewlett-Packard Labs*

## Abstract

In this paper, we rethink how an OS supports virtual memory. Classical VM is an opaque abstraction of RAM, backed by demand paging. However, most systems today (from phones to data-centers) do not page, and indeed may require the performance benefits of non-paged physical memory, precise NUMA allocation, etc. Moreover, MMU hardware is now useful for other purposes, such as detecting page access or providing large page translation. Accordingly, the venerable VM abstraction in OSes like Windows and Linux has acquired a plethora of extra APIs to poke at the policy behind the illusion of a virtual address space.

Instead, we present Cichlid, a memory system which inverts this model. Applications explicitly manage their physical RAM of different types, and directly (though safely) program the translation hardware. Cichlid is implemented in Barrelfish, requires no virtualization support, and outperforms VMM-based approaches for all but the smallest working sets. We show that Cichlid enables use-cases for virtual memory not possible in Linux today, and other use-cases are simple to program and significantly faster.

## 1   Introduction

We argue that applications for modern machines should manage physical RAM explicitly and directly program MMUs according to their needs, rather than manipulating such hardware implicitly through a virtual address abstraction as in Linux. We show that explicit primitives for managing physical memory and the MMU deliver comparable or better application performance, greater functionality, and a simpler and orthogonal interface that avoids the feature interaction and performance anomalies seen in Linux.

Traditional virtual memory (VM) systems present a conceptually simple view of memory to the application programmer: a single, uniform virtual address space which the OS transparently backs with physical memory. In its pure form, applications never see page faults, RAM allocation, address translation, TLB misses, etc.

This simplicity has a price. VM is an illusion — one can exhaust physical memory, resulting in thrashing, or the OS killing the application. Moreover, performance is unpredictable. VM hardware is complex, with multiple caches, TLBs, page sizes, NUMA nodes, etc.

For applications like databases the performance gains from closely managing the MMU mappings and locations of physical pages on memory controllers are as important to the end user as the functional correctness of the program [39,56]. Consequently, the once-simple VM abstraction in systems such as Linux has become steadily more complex, as application developers demand more control over the mapping hardware, by piercing the VM abstraction with features like transparent huge pages, NUMA allocation, pinned mappings, etc. In Section 2, we discuss the complexity, redundancy, and feature interaction in the formerly simple VM interface.

In response, we investigate the consequences of turning the VM system inside-out: applications (1) directly manage physical RAM, and (2) directly (but safely) program MMUs to build the environment in which they operate. Our contribution is a comprehensive design which achieves these goals, allows the full range of use-cases for memory system hardware, and which performs well.

Cichlid[1], a new memory management system built in the Barrelfish research OS, adopts a radically inverted view of memory management compared with a traditional system like Linux. Cichlid processes still run inside a virtual address space (the MMU is enabled) but this address space is securely constructed by the application itself with the help of a library which exposes the full functionality of the MMU. Above this, all the functionality of a traditional OS memory system is provided.

---

[1]Pronounced ˈsɪklɪd; see https://en.wikipedia.org/wiki/Cichlid.

Application-level management of the virtual address space is not a new idea; we review its history in Section 5. Cichlid itself is an extension of the original Barrelfish physical memory management system described in Baumann et al. [8], which itself was based on seL4 [32].

The contributions of Cichlid over these prior systems are:

- A comprehensive implementation of application-level memory management for modern hardware capable of supporting applications which exploit its features. We extend the Barrelfish model to support safe user construction of page tables, arbitrary superpage mapping, demand paging, and fast access to page status information without needing virtualization hardware.

- A detailed performance evaluation of Cichlid comparing it with a variety of techniques provided by, and different configurations of, a modern Linux kernel, showing that useful performance gains are achieved while greatly simplifying the interface.

In the next section of this paper we first review the various memory management features in Linux as an example of the traditional Unix-based approach. In Section 3 we then present Cichlid, and evaluate its performance in Section 4. Section 5 discusses the prior work on explicit physical memory management, and Section 6 summarizes the contribution and future work.

## 2 Background: the Linux VM system

We now discuss traditional VM systems, as context for Cichlid. We focus on Unix-like systems and Linux in particular as representative of mainstream approaches and the problems they exhibit. Later, in Section 5 we discuss prior systems which have adopted a different approach, some of which have strongly influenced Cichlid.

### 2.1 Traditional Unix

Unix was designed when RAM was scarce, and demand paging essential to system operation. Virtual memory is fully decoupled from backing storage via paging. Each process sees a uniform virtual address space. All memory is paged to disk by a single system-wide policy. The basic virtual memory primitive visible to software is `fork()`, which creates a complete copy of the virtual address space. Modern `fork()` is highly optimized (e.g. using copy-on-write).

Today, RAM is often plentiful, MMUs are sophisticated and featureful devices (e.g. supporting superpages), and the memory system is complex, with multiple controllers

and set-associative caches (e.g. which can be exploited with page coloring).

Workloads have also changed. High-performance multicore code pays careful attention to locality and memory controller bandwidth. Pinning pages is a common operation for performance and correctness reasons, and personal devices like phones are often designed to not page at all.

Instead, the MMU is used for purposes aside from paging. In addition to protection, remapping, and sharing of physical memory, MMUs are used to interpose on main memory (e.g. for copy-on-write, or virtualization) or otherwise record access (such as the use of "dirty" bits in garbage collection).

### 2.2 Modern Linux

The need to exploit the memory system fully is evident from the range of features added to Linux over the years to "poke through" the basic Unix virtual address abstraction.

The most basic of these creates additional "shared-memory objects" in a process' address space, which may or may not be actually shared. Such segments are referred to by file descriptors and can either be backed by files or "anonymous". The basic operation for mapping such an object is `mmap()`, which in addition to protection information accepts around 16 different flags specifying whether the mapping is shared, at a fixed address, contains pre-zeroed memory, etc. We describe basic usage of `mmap()` and related calls in Section 2.3; above this are a number of extensions.

**Large pages:** Modern MMUs support mappings at a coarser granularity than individual pages, typically by terminating a multi-level page table walk early. For example, `x86_64` supports 2 MB and 1 GB *superpages* as well as 4 kB pages, and for simplicity we assume this architecture in the discussion that follows (others are similar).

Linux support for superpage mappings is somewhat complex. Firstly, mappings can be created for large (2 MB) or huge (1 GB) pages via a file system, `hugetlbfs` [40, 58] either directly or through `libhugetlbfs` [41]. For each supported superpage size, a command-line argument tells the kernel to allocate a fixed pool of superpages at boot-time. This pool can be dynamically resized by an administrator. Shrinking a pool deallocates superpages from applications using a hard-wired balancing policy. In addition, one superpage size is defined as a system-wide default which will be used for allocation if not explicitly specified otherwise.

Once an administrator has set up the page pools, users can be authorized to create memory segments with superpage mappings, either by mapping files created in the `hugetlbfs` file system, or mapping anonymous segments

with appropriate flags. Superpages may not be demand-paged [5].

The complexity of configuring different memory pools in Linux at boot has led to an alternative, *transparent huge pages* (THP) [27, 59]. When configured, the kernel allocates large pages on page faults if possible according to a single, system-wide policy, while a low-priority kernel thread scans pages for opportunities to use large pages through defragmentation. Demand-paging is allowed by first splitting the superpage into 4 kB pages [5]. A typical modern `x86_64` kernel is configured for transparent support of 2 MB pages, but not 1 GB pages. Alternatively, an administrator can disable system-wide THP at boot or by writing to sysfs and programs can enable it on a per-region basis at runtime using `madvise()`.

**NUMA:** The `mbind()` system call sets a NUMA policy for a specific virtual memory region. A policy consists of a set of NUMA nodes and a mode: *bind* to restrict allocation to the given nodes; *preferred* to prefer those nodes, but fall back to others; *interleave* allocations across the nodes, and *default* to lazily allocate backing memory on the local node of the first thread to touch the virtual addresses. This "first touch" policy has proved problematic for performance [29].

`libNUMA` provides an additional `numa_alloc_onnode()` call to allocate anonymous memory on a specific node with `mmap()` and `mbind()`. Linux can move pages between nodes: `migrate_pages()` attempts to move all pages of a process that reside on a set of given nodes to another set of nodes, while `move_pages()` moves a set of pages (specified as an array of virtual addresses) to a set of nodes. Note that policy is expressed in terms of virtual, not physical, memory.

There are also attempts [19–22, 25, 29] to deal with NUMA performance issues transparently in the kernel, by migrating threads closer to the nodes containing memory they frequently access, or conversely migrating pages to threads' NUMA nodes, based on periodically revoking access to pages and tracking usage with soft page faults. A good generic policy, however, may be impossible; highly performance-dependent applications currently implement custom NUMA policies by modifying the OS [29].

**User-space faults:** Linux signals can be used to reflect page faults to the application. GNU libsigsegv [43] provides a portable interface for handling page faults: a user fault handler is called with the faulting virtual address and must then be able to distinguish the type of fault, and possibly map new pages to the faulting address. When used with system calls such as `mprotect()` and `madvise()`, this enables basic user-space page management. The current limitations of this approach (both in performance and flexibility) have led to a proposed facility for user-space demand paging [23, 26].
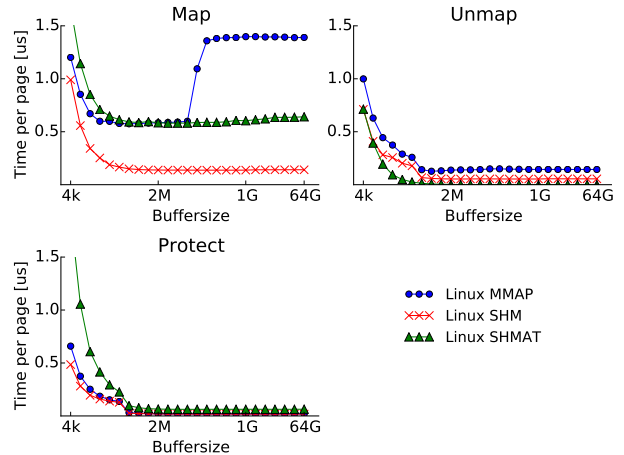


Figure 1: Managing memory on Linux (`4.2.0`)

## 2.3 Discussion

Based on the simple Unix virtual address space, the Linux VM system has evolved in response to new demands by accreting new features and functionality. This has succeeded up to a point, but has resulted in a number of problems.

The first is **mechanism redundancy**: there are multiple mechanisms available to users with different performance characteristics. For example, Figure 1 shows the performance of three different Linux facilities for creating, destroying, and changing "anonymous mappings": regions of virtual address space backed by RAM but not corresponding to a file. These measurements were obtained using the machine in Table 1 using 4k pages throughout.

MMAP uses an `mmap()` call with `MAP_POPULATE` and `MAP_ANONYMOUS` to map and unmap regions, and `mprotect()` for protection. This forces the kernel to zero pages being mapped, dominating execution time. Avoiding this behavior, even when safe, requires kernel reconfiguration at build time – a global policy aimed at embedded systems.

SHM creates a shared memory object with `shm_open()` and passes it to `mmap()` and `mprotect()`. In this case, `mmap()` will not zero the memory. Unmapping is also faster since memory is not immediately reclaimed. The object can be shared with other processes, but (unlike MMAP mappings) cannot use large pages.

SHMAT attaches a shared segment with `shmat()`, and *does* allow large pages if the process has the `CAP_IPC_LOCK` capability. Internally, the mechanism is similar to `mmap()`, with system-wide limits on the number and size of segments.

For buffers up to 2 MB, the cost per page decreases with size for all operations due to amortization of the system call overhead. Afterwards, the time stays constant

| CPU | Intel Xeon E5-2670 v2 (Ivy Bridge) |
|---|---|
| #nodes / #sockets / #cores | 2 / 2 / 20 @ 2.5 GHz |
| L1 / L2 cache | 32 kB / 256 kB (per core) |
| L3 size | 25 MB (shared) |
| dTLB (4 kB pages) | 64 entries (4-way) |
| dTLB (2 MB pages) | 32 entries (4-way) |
| dTLB (1 GB pages) | 4 entries (4-way) |
| L2 TLB (4K) | 512 entries (4-way) |
| RAM | 256 GB (128 GB per node) |
| Linux kernel | v.4.2.0 (Ubuntu 15.10) |

Table 1: Test bed specifications. [49]

| 4.2.0 | 4.2.0 (Ubuntu 15.10) | No large page support |
|---|---|---|
| 4.2.0-tlbfs | 4.2.0 (Ubuntu 15.10) | hugetlbfs enabled |
| 4.2.0-thp | 4.2.0 (Ubuntu 15.10) | Transparent huge pages enabled |
| 3.16 | 3.16 | Stock 3.16 kernel |
| 3.16-dune | 3.16 | Linux 3.16 with Dune |

Table 2: Tested Linux configurations

except for MMAP map operations.

libhugetlbfs provides get_hugepage_region and get_huge_pages calls to directly allocate superpage-backed memory using a malloc-style interface. The actual page size cannot be specified and depends on a system-wide default; 4 kB pages may be used transparently unless the GHR_STRICT flag is set. By default, hugetlbfs prefaults pages.

The high-level observation is: *No single Linux API is always optimal, even for very simple VM operations.*

A second problem is **policy inflexibility**. While the appropriate policy for many memory management operations such as page replacement, NUMA allocation or handling of superpages depend strongly on individual application's workloads. In Linux, however, they usually either apply system-wide, require administrator configuration (often at boot), must be enabled at compile time, or a combination of them.

For example, supporting two superpage sizes in hugetlbfs requires two different, pre-allocated pools of physical memory, each assigned to a different file system, precluding a dynamic algorithm that could adapt to changing workloads.

In addition to the added complexity in the kernel [24], the system-wide policies in *transparent* superpage support have led to a variety of performance issues: Oracle DB has suffered from I/O performance degradation when reading large extents from disk [5, 17]. Redis incurs unexpected latency spikes using THP due to copy-on-write overhead for large pages, since the application periodically uses fork() to persist database snapshots [65]. The jemalloc memory allocator experiences performance anomalies due to its use of madvise to release small regions of memory inside of bigger chunks which have been transparently backed by large pages — the resulting holes preventing later merging of the region back into a large page [37].

These issues are not minor implementation bugs, but arise from the philosophy that memory system complexity should be hidden from applications, and resource allocation policies should be handled transparently by the kernel.

The third class of problem is **feature interaction**. We have seen how superpages cannot be demand paged (even though modern SSDs can transfer 2MB pages with low latency). Another example is the complex and subtle interaction between kernel-wide policies for NUMA allocation with superpage support [58]. At one level, this shows up in the inability to control initial superpage allocation at boot time (superpages are always balanced over all NUMA nodes). Worse, Gaud et al. [38] show that treating large pages and NUMA separately does not work well: large pages hurt the performance of parallel applications on NUMA machines because *hot pages* are more likely, and larger, and *false page sharing* makes replication or migration less effective. Accordingly, the Carrefour [29] system modifies the kernel's NUMA-aware page placement to realize its performance gains.

Collectively, these issues motivate investigating alternative approaches. As memory hardware diversifies in the future, memory management policies will become increasingly complicated. We note that none of the Linux memory APIs actually deal with *physical* memory directly, but instead select from a limited number of complex, in-kernel policies for backing traditional *virtual* memory.

In contrast, therefore, Cichlid safely exposes to programs and runtime systems both physical memory and translation hardware, and allows libraries to build familiar virtual memory abstractions above this.

## 3 Design

We now describe the design of Cichlid, and how it is implemented over the basic memory functionality of Barrelfish. While Cichlid allows great flexibility in arranging an address space, it nevertheless ensures the following key safety property: *no Cichlid process can issue read or write instructions for any area of physical memory for which it does not have explicit access rights.*

Subject to this requirement, Cichlid also provides the following completeness property: *a Cichlid process can create any address space layout permitted by the MMU for which it has sufficient resources.* In other words, Cichlid itself poses no restriction on how the memory hardware can be used.

There are three main challenges in the implementa-

tion that Cichlid must address: Firstly, it must securely name and authorize access to, and control over, regions of physical memory. Cichlid achieves this using *partitioned capabilities*. Secondly, it must allow safe control of hardware data structures (such as page tables) by application programs. This, is achieved by considerably extending the set of memory types supported by the capability system in Barrelfish (and seL4) for Cichlid to use. Finally, Cichlid must give applications direct access to information provided by the MMU (such as access and write-tracking bits in the page tables). Unlike prior approaches which rely on virtualization technology, Cichlid allows direct read-only access to page table entries; we explain below why this is safe.

Cichlid has three main components: First, the kernel provides capability invocations that allow application processes to install, modify and remove page table entries and query for the base address and size of physical regions. Second, the kernel exception handler redirects any exceptions generated by the MMU to the application process that caused the exception. Thirdly, a runtime library provides to applications an abstraction layer over the capability system which exposes a simple, but expressive API for managing page tables.

## 3.1   Physical memory allocation

Cichlid applications directly allocate regions of physical memory and pass around authorization for these regions in the form of capabilities. Regions can be mapped into a virtual address space by changing a page table, or used for other purposes such as holding page tables themselves.

Cichlid extends the Barrelfish capability design, itself inspired by seL4 [30, 33, 53]. All physical regions are represented by capabilities, which also confer a particular *memory type*. For example, the integrity of the capability system itself is ensured by storing capability representations in memory regions of type `CNode`, which can never be directly written by user-space programs. Instead, a region must be of type `Frame` to be mapped writable into a virtual address space. Holding both `Frame` and `CNode` capabilities to the same region would enable a process to forge new capabilities by directly manipulating their bit representations, and so is forbidden. Such a situation is prevented by having a kernel enforced type hierarchy for capabilities.

Capabilities to memory regions can be split and *retyped* according to a set of rules. At system start-up, all memory is initially of type `Untyped`, and physical memory is allocated to processes by splitting the initial untyped region. Retyping and other operations on capabilities is performed by system calls to the kernel.

seL4 capabilities are motivated by the desire to prove correctness properties of the seL4 kernel, in particular,

the property that no system call can fail due to lack of memory. Hence, seL4 and Barrelfish perform no dynamic memory allocation in the kernel, instead memory for all dynamic kernel data structures is allocated by user-space programs and retyped appropriately, such as to a kernel thread control block or a `CNode`, for example.

Capabilities are attractive since they export physical memory to applications in a safe manner: application may not arbitrarily use physical memory; they must instead "own" the corresponding capability. Furthermore, capabilities can be passed between applications. Finally, capabilities have some characteristics of objects: each capability type has a set of *operations* which can be invoked on it by a system call.

In Barrelfish, seL4, and Cichlid, the kernel enforces safety using two types of meta-data: a *derivation database* and a per-processes *capability space*. All capability objects managed by a kernel are organized in a capability derivation tree. This tree enables efficient queries for descendants (of retype and split operations) and copies. These queries are used to prevent retype races on separate copies of a capability that might compromise the system.

User processes refer to capabilities and invoke operations on them using opaque handles. Each process has its own capability address space, which is explicitly maintained via a radix tree in the kernel which functions as a *guarded page table*. The nodes of the tree are also capabilities (retyped from RAM capabilities) and are allocated by the application.

The root of the radix tree for each process is stored in the process control block. When a process invokes a capability operation it passes to the kernel the capability handle with the invocation arguments. To perform the operation, the kernel traverses the process' capability space to locate the capability corresponding to the handle and authorizes the invocation.

Cichlid builds on the basic Barrelfish capability mechanisms to allow explicit allocation of different kinds of memory. A memory region has architectural attributes such as the memory controller it resides on, whether it is on an external co-processor like a GPGPU or Intel Xeon Phi, whether it is persistent, etc. Applications explicitly acquire memory with particular attributes by requesting a capability from an appropriate memory allocator process, of which there are many. Furthermore, less explicit "best effort" policies can be layered on top by implementing further virtual allocators which can, for example, steal RAM from nearby controllers if local memory is scarce.

## 3.2   Securely building page tables

Page tables are hardware specific, and at the lowest level, Cichlid's interface (like seL4 and Barrelfish) reflects the actual hardware. Applications may use this interface di-

rectly, or a high-level API with common abstractions for different MMUs, to safely build page tables, exchange page tables on a core, and install mappings for any physical memory regions for which the application is authorized. The choice of virtual memory layout, and its representation in page tables, is fully controlled by the application. Cores can share sub-page-tables between different page-table hierarchies to alias a region of memory at a different address or to share memory between different cores as in Corey [16].

Cichlid adds support for multiple page sizes (2 MB and 1 GB superpages in x86_64, and 16 MB, 1 MB, and 64 kB pages in ARMv7-a [4]) to the Barrelfish memory management system [8]. Cichlid decouples the physical memory allocation from programming the MMU. Therefore the API allows for a clean way to explicitly select the page size for individual mappings, map pages from a mixture of different page sizes, and change the virtual page sizes for mappings of contiguous physical memory regions all directly from the applications itself instead of relying on the kernel to implement the correct policy for all cases.

To do this, Cichlid extends the Barrelfish memory system (and that of seL4) by introducing a new capability type for every level of page table for every architecture supported by the OS. This is facilitated by the *Hamlet* domain-specific language for specifying capability types [28].

For example, for an MMU in x86_64 long-mode there are four different types of page table capability, corresponding to the 4 levels of a 64-bit x86 page table (PML4, PDPT, PD, and PT). A PT (last-level page table) capability can only refer to a 4k page-aligned region of RAM and has a map operation which takes an additional capability plus an entry number as arguments. This capability in turn must be of type Frame and refer to another 4k page. The operation installs the appropriate page table entry in the PT to map the specified frame. The kernel imposes no policy on this mapping, other than restricting the type and size of capabilities.

Similarly, a map on a PD (a 2nd-level "page directory") capability only accepts a capability argument which is of size 4 kB and type PT, *or* of type Frame and size 2 MB (signifying a large page mapping).

A small set of rules therefore captures all possible valid and authorized page table operations for a process, while excluding any that would violate the safety property. Moreover, checking these rules is fast and is partly responsible for Cichlid's superior performance described in Section 4.2. This type system allows user-space Cichlid programs to construct flexible page tables while enforcing the safety property stated at the start of this section.

Cichlid's full kernel interface contains the following capability invocations: identify, map, unmap, modify_flags (protect), and clear_dirty_bits.

Memory regions represented by capabilities and associated rights allow user-level applications to safely construct page tables; they allocate physical memory regions and retype them to hold a page table and install the entries as needed.

Typed capabilities ensure a process cannot successfully map a physical region for which it does not have authorization. The process of mapping itself is still a privileged operation handled by the kernel, but the kernel must only validate the references and capability types before installing the mapping. Safety is guaranteed based on the type system: page tables have a specific type which cannot be mapped writable.

Care must be taken in Cichlid to handle capability revocation. In particular, when a Frame capability is revoked, all page table entries for that frame must be quickly identified and removed. Cichlid handles this by requiring each instance of a Frame capability to correspond to at most one hardware page table entry. To map a frame into multiple page tables, or at multiple locations in the same page table, the program must explicitly create copies of the capability.

As described so far, each operation requires a separate system call. Cichlid optimizes this in a straightforward way by allowing batching of requests, amortizing system call cost for large region operations. The map, unmap, and modify_flags operations all take multiple consecutive entries for a given page table as arguments.

In Section 4.3 we confirm existing work on the effect of page size on performance of particular workloads, and in Section 4.4 we show that the choice of the page size is highly dynamic and depends on the program's configuration such as the number of threads and where memory is allocated.

In contrast, having the OS transparently select a page size is an old idea [60] and is the default in many Linux distributions today, but finding a policy that satisfies a diverse set of different workloads is difficult in practice and leads to inherent complexity with questionable performance benefits [17, 38, 42, 65].

### 3.3 Page faults and access to status bits

Cichlid uses the existing Barrelfish functionality for reflecting VM-related processor exceptions back to the faulting process, as in Nemesis [44] and K42 [55]. This incurs lower kernel overhead than classical VM and allows the application to implement its own paging policies. In Sections 4.1 and 4.2 we show that Cichlid's trap latency to user space is considerably lower than in Linux.

Cichlid extends Barrelfish to allow page-traps to be eliminated for some use-cases when the MMU maintains page access information in the page table entries. While Dune [9] uses nested paging hardware to present "dirty"

and "accessed" bits in an `x86_6 4` page table to a user space program, Cichlid achieves this *without* hardware support for virtualization.

We extend the kernel's mapping rules in Section 3.2 to allow page tables themselves to be mapped read-only into a process' address space. Essentially, this boils down to allowing a 4 kB capability of type `PML4`, `PDPT`, `PD`, or `PT` to be mapped in an entry in a PT instead of a `Frame` capability, with the added restriction that the mapping must be read-only.

This allows applications (or libraries) to read "dirty" and "accessed" bits directly from page table entries without trapping to the kernel. Setting or clearing these bits remains a privileged operation which can only be performed by a kernel invocation passing the capability for the page table.

Note that this functionality remains safe under the capability system: an application can only access the mappings it has installed itself (or for which it holds a valid capability), and cannot subvert them.

In Section 4.5 we demonstrate the benefits of this approach for a garbage collector. Cichlid's demand-paging functionality for `x86_64` also uses mapped dirty bits to determine if a frame's contents should be paged-out before reusing the frame.

Since Cichlid doesn't need hardware virtualization support, such hardware, if present, can be used for virtualization. Cichlid can work both inside a virtual machine, or as a better memory management system for a low-level hypervisor.

Moreover, nested paging has a performance cost for large working sets, since TLB misses can be twice as expensive. In Section 4.6 we show that for small working sets (below 16 MB for our hardware) a Dune-like approach outperforms Cichlid due to lower overhead in clearing page table bits, but for medium-to-large working sets Cichlid's lower TLB miss latency improves performance.

The Cichlid and Dune approaches are complementary, and a natural extension to Cichlid (not pursued here) would allow applications access to both the physical (machine) page tables *and* nested page tables if the workload can exploit them.

### 3.4 Runtime library

Cichlid provides a number of APIs above the capability invocations discussed above.

The first layer of Cichlid's user-space wraps the invocations to create, modify or remove single mappings, and keep track of the application's virtual address space layout.

While application programmers can build directly on this, the Cichlid library provides higher-level abstractions

based on the concepts of *virtual regions* (contiguous sets of virtual addresses), and *memory objects* that can be used to back one or more virtual regions and can themselves be comprised of one or more physical regions.

This layer is important to Cichlid's usability. Manually invoking operations on capabilities to manage the virtual address space can be cumbersome; take the example of a common operation such as mapping an arbitrarily-sized region of physical memory $R$ with physical base address $P$ and size $S$ bytes, $R = (P, S)$, at an arbitrary virtual base address $V$. The number of invocations needed to create this simple mapping varies based on $V$, $S$, and the desired properties of the mapping (such as page size), as well as the state of the application's virtual address space before the operation. In particular, installing a mapping can potentially entail creating multiple levels of page table in addition to installing a page table entry. The library encapsulates the code to do this on demand, as well as batching operations up to amortize system call overhead.

Finally, the library also provides traditional interfaces such as `sbrk()` and `malloc()` for areas of memory where performance is not critical. To simplify start-up, programs running over Cichlid start up with a limited, conventional virtual address space with key segments (text, data, bss) backed with RAM, though this address space is, itself, constructed by the process' parent using Cichlid (rather than the kernel).

In addition, the Cichlid library provides demand paging to disk as in Nemesis [44], but not by default: many time-sensitive applications rely on *not* paging for correctness, small machines such as phones typically do not page anyway, and and the growth of non-volatile main memory [48, 62] may make demand-paging obsolete. Unlike the Linux VM system, demand paging is orthogonal to page size: the Cichlid library can demand-page superpages provided the application has sufficient frames and disk space. Furthermore, the application is aware of the number of backing frames and can add or remove frames explicitly at runtime if required.

The library shows that building a classic VM abstraction over Cichlid is straightforward, but the reverse is not the case.

## 4 Evaluation

We evaluate Cichlid by first demonstrating that primitive operations have performance as good as, or better than those of Linux, and then showing that Cichlid's flexible interface allows application programmers to usefully optimize their systems.

All Linux results, other than those for Dune (Section 4.5), are for version 4.2.0, as shipped with Ubuntu 15.10, with three large-page setups: none, `hugetlbfs`, and transparent huge pages. As the Dune patches (git revision
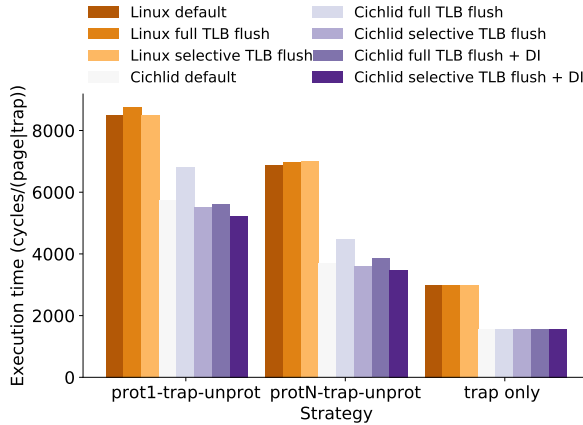
Figure 2: Appel-Li benchmark. (Linux `4.2.0`)

6c12ba0) require a version 3 kernel, these benchmarks use kernel version 3.16 instead. These configurations are summarized in Table 2. Thread and memory pinning was done using `numactl` and `taskctl`. Performance numbers for Linux are always the best among all tested configurations.

## 4.1 Appel and Li benchmark

The Appel and Li benchmark [3] tests operations relevant to garbage collection and other non-paging tasks. This benchmark is compiled with flags `-O2 -DNDEBUG`, and summarized in Figure 2.

We compare Linux and Cichlid with three different TLB flush modes: 1) Full: Invalidate the whole TLB (writing `cr3` on `x86`) every time, 2) Selective: Only invalidate those entries relevant to the previous operation (using the `invlpg` instruction), and 3) System default: Cichlid, by default, does a full flush only for more than one page. Linux's default behavior depends on kernel version. The version tested (4.2.0) does a selective flush for up to 33 pages, and full a flush otherwise [45]. We vary this value to change Linux's flush mode. The working set here is less than 2 MB, and thus large pages have no effect and are disabled.

Cichlid is consistently faster than Linux here.

For multi-page protect-trap-unprotect (*protN-trap-unprot*), Cichlid is 46% faster than Linux. For both systems, the default adaptive behavior is as good as, or better than, selective flushing. The Cichlid *+DI* results use the kernel primitives directly, to isolate the cost of user-space accounting, which is around 5%.

## 4.2 Memory operation microbenchmarks

We extend the Appel and Li benchmarks, to establish how the primitive operations scale for large address spaces, using buffers up to 64 GB. We *map*, *protect* and *unmap* the entire buffer, and time each operation separately. We compare Cichlid to the best Linux method for each page size established in § 2.3. On Cichlid we use the high-level interfaces on a previously allocated frame, for similar semantics to shared memory objects in Linux. The experiments were conducted on a 2x10 Intel Xeon E5 v2. Figure 3 shows execution time per page.

**Map:** Cichlid per-page performance is highly predictable, regardless of page size. Since all information needed is presented to each a system call, the kernel does very little. On Linux we use `shm_open` for 4k pages and `shmat` for others. Linux needs to consult the shared segment descriptor and validate it. This results in a general performance improvement for Cichlid over Linux up to 15x for 4 kB pages or 93x for large pages, once some upfront overhead is amortized.

**Protect:** These are in line with the Appel and Li benchmarks: Cichlid outperforms Linux's `mprotect()` on an `mmap`'ed region in all configurations except for small buffers of 4 kB pages. For large buffers, the differences between Cichlid and Linux are up to 4x (4 kB pages) or 8x (huge pages).

**Unmap:** Doing an unmap in Cichlid is expensive: the relevant page table capability must be looked up to invoke it and the mapped `Frame` capability needs to be marked unmapped. Linux `shmdt`, however, simply detaches the segment from the process but doesn't destroy it. Cichlid could be modified to directly invoke the page table, and thereby match the performance of Linux.

Cichlid memory operations are competitive: capabilities and fast traps allows an efficient virtual memory interface. Even when multiple page table levels are changed, Cichlid usually outperforms Linux on most cases, despite requiring several system calls.

## 4.3 Random accesses benchmark (GUPS)

Many HPC workloads have a random memory access pattern, and spend up to 50% of their time in TLB misses [67]. Using the RandomAccess benchmark [54] from the HPC Challenge [68] suite, we demonstrate that carefully user-selected page sizes, as enabled by Cichlid, have a dramatic performance effect.

We measure update rate (Giga updates per second, or GUPS) for read-modify-write on an array of 64-bit integers, using a single thread. We measure working sets up to 32 GB, which exceeds TLB coverage for all page sizes. Linux configuration is `4.2.0-tlbfs`, with pages allocated from the local NUMA node. If run with transpar-
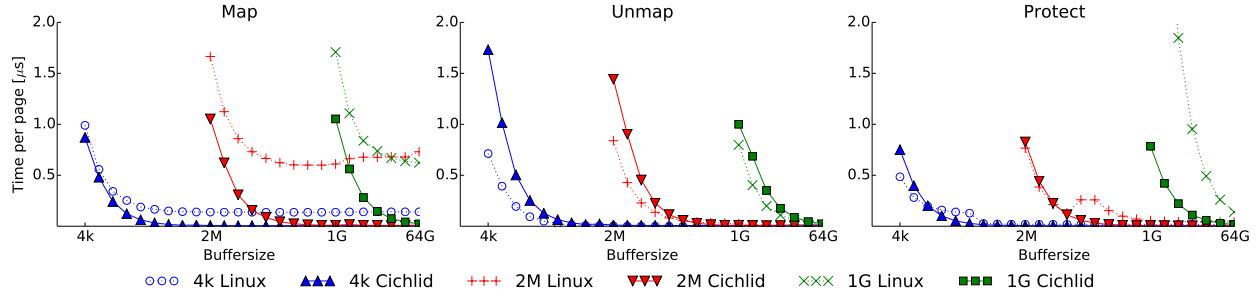
Figure 3: Comparison of memory operations on Cichlid and Linux using `shmat`, `mprotect` and `shmdt`. (Linux `4.2.0-tlbfs`)

| Page Size | Cichlid | | Linux | |
|---|---|---|---|---|
| | GUPS | Time | GUPS | Time |
| 4k | 0.0122 | 1397s | 0.0121 | 1414s |
| 2M | 0.0408 | 420s | 0.0408 | 421s |
| 1G | 0.0659 | 260s | 0.0658 | 261s |

Table 3: GUPS as a function of page size, 32 GB table.



Figure 4: GUPS as a function of table size, normalized.



Figure 5: GUPS variance. `4.2.0-tlbfs`, 2 MB pages.

ent huge pages instead, the system always selects 2 MB pages, and achieves lower performance.

Figure 4 shows the results on Cichlid, normalized to 1 GB pages. Performance drops once we exceed TLB coverage: at 2 MB for 4 kB pages, and at 128 MB for 2 MB pages. The apparent improvement at 32 MB is due to exhausting the L3 cache, which slows all three equally, bringing the normalized results together. Large pages not only increase TLB coverage, but cause fewer table walk steps to service a TLB miss. Page-structure caches would reduce the number of memory accesses even further but are rather small [6, 12] in size. Cichlid and Linux perform identically in the test, as Table 3 shows. These results support previous findings on TLB overhead [7, 67], and emphasize the importance for applications being able to
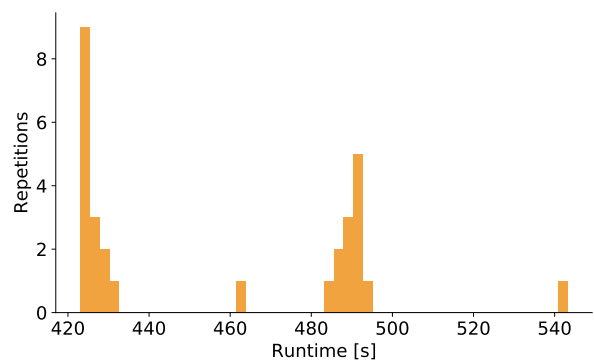
select the correct page size for their workload.

On Linux, even with NUMA-local memory, high scheduling priority, and no frequency scaling or power management, there is a significant variance between benchmark runs, evidenced by the multimodal distribution in Figure 5. This occurs for both `hugetlbfs` and transparent huge pages, and is probably due to variations in memory allocation, although we have been unable to isolate the precise cause. This variance is completely absent under Cichlid even when truly randomizing paging layout and access patterns, demonstrating again the benefit of predictable application-driven allocation.

## 4.4 Mixed page sizes

Previous work [38] has shown that while large pages can be beneficial on NUMA systems, they can also hurt performance. Things are even more complicated when there are more page sizes (e.g., 4 kB, 2 MB, 1 GB for `x86_64`). Furthermore, modern machines often have a distinct TLB for each page size, suggesting that using a mix of page sizes increases TLB coverage.

Kaestle *et al.* [50] showed that distribution and repli-

| CPU | AMD Opteron 6378 |
| micro architecture | Piledriver |
| #nodes / #sockets / #cores | 8 / 4 / 32 @ 2.4 GHz |
| L1 / L2 cache size | 16 kB / 2 MB per core |
| L3 cache size | 12 MB per socket |
| dTLB (4 kB pages) | 64 entries, fully |
| dTLB (2/4 MB pages) | 64 entries, fully |
| dTLB (1 GB pages) | 64 entries, fully |
| L2 TLB (4 kB pages) | 1024 entries, 8 way |
| L2 TLB (2/4 MB pages) | 1024 entries, 8 way |
| L2 TLB (1 GB pages) | 1024 entries, 8 way |
| RAM | 512 GB (64 GB per node) |

Table 4: Specification of machine used in §4.4

| page size | array configuration | | |
| --- | --- | --- | --- |
| | T=1 | T=32 (dist) | T=32 (repl + dist) |
| 4 kB | 597.91 | **51.32** | 34.43 |
| 2 MB | 414.80 | 58.09 | **28.87** |
| 1 GB | **395.64** | 265.94 | 128.77 |

Table 5: PageRank runtime (seconds) depending on page size and PageRank configuration (repl = replication, dist = distribution, T is the number of threads). Highlighted are best numbers for each configuration. Standard error is very small.

cation of data mitigates congestion on interconnects and balances memory controller load, by extending Green-Marl [47], a high-level domain-specific language for graph analytics, to automatically apply these techniques per region, using patterns extracted by the compiler. This gave a two-fold speedup of already tuned parallel programs.

Large pages interact with the NUMA techniques described above, by changing the granularity at which they can be applied to data structures that are contiguous in virtual memory. The granularity of NUMA distribution, for example, is the page size. Hence, the smaller the page size the more slack the run-time has to distribute data across NUMA nodes. Bigger page sizes also make memory allocation more restrictive: The starting address when allocating memory must be a multiple of the page size. Bigger page sizes can increase fragmentation and increases the chance of conflicts in caches and TLB.

In Cichlid, programs map their own memory, and all combinations of page sizes are supported. Furthermore, no complex setup of page allocations and kernel configurations are required.

Table 5 shows the effect of the page size on application performance using Shoal's Green-Marl PageRank [50]. NUMA effects are minimal on the 2-socket machine we are using in other experiments, so for this experiment we use the machine in Table 4 and note that AMD's SMT threads (CMT) are disabled in our experiments.

We evaluate two configurations: First, single-threaded (T=1). In this case replication does not make sense as all accesses are local, and distribution is unnecessary as a single thread cannot saturate the memory controller — indeed, an increase in remote memory access would likely reduce performance. In this case, an isolated application, bigger pages are always better.

Next, we run on all cores and explore the impact of replication and distribution on the choice of page sizes. 1 GB pages clearly harm performance as distribution is impossible or too coarse-grained. We only break even if

90% of the working set is replicated. However, the last 10% still cannot be distributed efficiently, which leads to worse performance.

It is clear that the right page size is highly dynamic and depends on workload and application characteristics. It is impractical to statically configure a system with pools (as in Linux) optimally for all programs, as the requirements are not known beforehand. Also, memory allocated to pools is not available for allocations with different page sizes. In contrast, Cichlid's simpler interface allows arbitrary use of page sizes and replication by the application without requiring *a priori* configuration of the OS.

## 4.5 Page status bits

The potential of using the MMU to improve garbage collection is known [3]. Out of many possible applications, we consider detecting page modifications; A feature used, for example, in the Boehm garbage collector [15] to avoid stopping the world. Only after tracing does the collector stop the world and perform a final trace that need only consider marked objects in dirty pages. This way, newly reachable objects are accounted for and not collected.

There are two ways to detect modified pages: The first is to make the pages read-only (e.g., via `mprotect()` or transparently by the kernel using soft-dirty PTEs [66]), and handle page faults in user-space or kernel-space. The handler sets a virtual dirty bit, and unprotects the page to allow the program to continue. The second approach uses hardware dirty bits, set when a page is updated. Some OSes (e.g., Linux) do not provide access to these bits. This is not just an interface issue. The bits are actively used by Linux to detect pages that need to be flushed to disk during page reclamation. Other OSes such as Solaris expose these dirty bits in a read-only manner via the `/proc` file-system. In this case, applications are required to perform a system call to read the bits, which, can lead to worse performance than using `mprotect()` [13].

In Cichlid, physical memory and page tables are directly visible to applications. Applications can map page tables read-only in their virtual address space. Only clearing the dirty bits requires a system call.
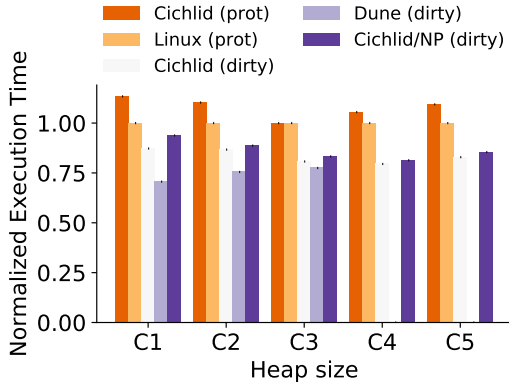
Figure 6: GCBench on Linux, Cichlid and Dune, normalized runtime to Linux. (Linux `3.16`, `3.16-dune`)

| Config | C1 | C2 | C3 | C4 | C5 |
|---|---|---|---|---|---|
| Runtime (s) | | | | | |
| Linux (prot) | 2.1 | 9.6 | 42 | 191 | 848 |
| Cichlid (prot) | 2.4 | 10.5 | 43 | 203 | 928 |
| Cichlid (dirty) | 1.9 | 8.3 | 34 | 153 | 692 |
| Dune (dirty) | 1.5 | 7.3 | 33 | – | – |
| Cichlid/NP (dirty) | 2.0 | 8.6 | 36 | 157 | 720 |
| Collections | | | | | |
| Linux (prot) | 251 | 336 | 381 | 428 | 448 |
| Cichlid (prot) | 245 | 335 | 393 | 432 | 442 |
| Cichlid (dirty) | 230 | 323 | 383 | 435 | 441 |
| Dune (dirty) | 318 | 367 | 403 | – | – |
| Cichlid/NP (dirty) | 233 | 325 | 381 | 434 | 443 |
| Heap size (MB) | | | | | |
| Linux (prot) | 139 | 411 | 1924 | 7972 | 24932 |
| Cichlid (prot) | 132 | 453 | 1413 | 6789 | 26821 |
| Cichlid (dirty) | 100 | 453 | 1477 | 5669 | 28132 |
| Dune (dirty) | 106 | 386 | 1579 | – | – |
| Cichlid/NP (dirty) | 100 | 453 | 1573 | 5541 | 28132 |

Table 6: GCBench reported total runtime, heap size and amount of collections.

Dune [9] provides this functionality through nested paging hardware, intended for virtualization, by running applications as a guest OS. Dune applications have direct access to the virtualized (nested) page tables. This approach avoids any system call overhead to reset the dirty bits, but depends on virtualization hardware and can lead to a performance penalty due to greater TLB usage [7, 11].

We use the Boehm garbage collector [15] and the GCBench microbenchmark [14]. GCBench tests the garbage collector by allocating and collecting binary trees of various sizes. We run this benchmark with the three described memory systems, Linux, Dune and Cichlid with five different configurations C1 to C5, which progressively increase the size of the allocated trees.

In Figure 6 we compare the runtime of each system. Cichlid implements all three mechanisms: protecting pages (Cichlid (prot)), hardware dirty bits (Cichlid (dirty)) in user-space and hardware dirty bits in guest ring 0 (Cichlid/NP (dirty)) (as does Dune). Our virtualization code is based on Arrakis [63].

Cichlid (prot) performs slightly worse than Linux (prot). This is consistent with Figure 3 where Linux performs better than Cichlid for protecting a single 4 kB page. We achieve better performance (between 13% (C2) and 19% (C4)) than Linux when we use hardware dirty bits, by avoiding traps when writing to pages. We still incur some overhead as we have to make a system call to reset the dirty bits on pages. Dune outperforms Cichlid (dirty) by up to 21% (C1), as direct access to the guest page tables enables resetting the dirty bits without having to make a system call. However, Cichlid manages to close the gap as the working set becomes larger, in which case Dune performance noticeably shows the overhead of nested paging. Unfortunately, we were unable to get Dune working with larger heap sizes on our hardware and thus have no numbers for Dune for configurations C4 and

C5.

On Linux, using transparent huge pages did not have a significant impact on performance and we report the Linux numbers with THP disabled. In a similar vein, we were unable to get Dune working with superpages, but we believe that having superpages might improve Dune performance for larger heap sizes (c.f. 4.3).

Cichlid/NP (dirty) runs GCBench in guest ring 0 and reads and clears dirty bits directly on the guest hardware page tables. The performance for Cichlid/NP is similar to Cichlid (dirty) and slower than Dune. However, this can be attributed to the fact that Cichlid/NP does not fully leverage the advantage of having direct access to the guest hardware page tables and still uses system calls to construct the address space.

Table 6 shows the total runtime, number of collections the GC did and the heap size used by the application. Ideally, the heap size should be identical for all systems since it is always possible to trade memory for better run time in a garbage collector. In practice this is very difficult to enforce especially across entirely different operating systems. For example Cichlid uses less memory (28%) for C4 compared to Linux (prot) but more memory (12%) for C5.

We conclude that with Cichlid we can safely expose MMU information to applications which in turn can benefit from it without relying on virtualization hardware features.
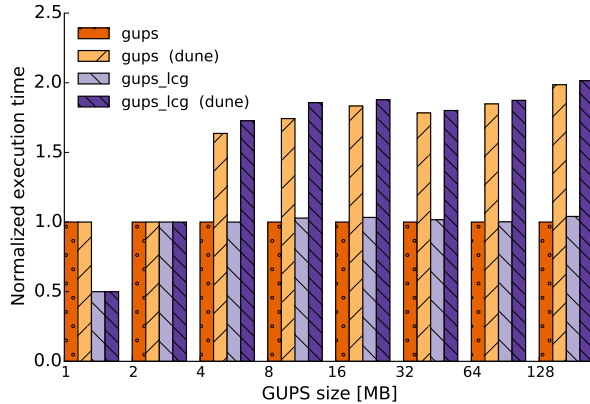
Figure 7: Comparison of the execution time of RandomAccess with and without nested paging for varying working set sizes, normalized to GUPS on native Linux. (Linux `3.16`, `3.16-dune`)

| Size | Linux | | Dune | |
| --- | --- | --- | --- | --- |
| | GUPS | GUPS LCG | GUPS | GUPS LCG |
| 1 | 2 | 1 | 2 | 1 |
| 2 | 3 | 3 | 3 | 3 |
| 4 | 11 | 11 | 18 | 19 |
| 8 | 35 | 36 | 61 | 65 |
| 16 | 90 | 93 | 165 | 169 |
| 32 | 236 | 240 | 421 | 425 |
| 64 | 594 | 595 | 1098 | 1113 |
| 128 | 1510 | 1571 | 2999 | 3043 |

Table 7: RandomAccess absolute execution times in milliseconds. (Linux `3.16`, `3.16-dune`)

## 4.6 Nested paging overhead

To illustrate the potential downside of nested paging, we revisit the HPC Challenge RandomAccess benchmark. Resolving a TLB miss with nested paging requires a 2D page table walk and up to 24 memory accesses [2] resulting in a much higher miss penalty, and the overhead of nested paging may end up outweighing the benefits of direct access to privileged hardware in guest ring zero. GUPS represents a worst-case scenario due to its lack of locality.

We conduct the same experiment as in section 4.3 on Dune [9] with a working set size ranging from 1 MB to 128 MB. Figure 7 and Table 7 show that for the smallest table sizes (1 MB and 2 MB) the performance of RandomAccess under Dune and Linux is comparable. Larger working set sizes exceed the TLB coverage and hence more TLB misses occur. This results in almost 2x higher runtime for RandomAccess in Dune than Linux. As for all comparisons with Dune, we disable transparent huge pages on Linux.

Running applications in guest ring zero as in Dune has pros and cons: on one hand, the application gets access to privileged hardware features, on the other hand, the performance may be degraded due to larger TLB miss costs for working sets which cannot be covered by the TLB.

## 4.7 Page coloring

The core principle of paged virtual memory is that virtual pages are backed by arbitrary physical pages. This can adversely affect application performance due to unnecessary conflict misses in the CPU caches and an increase in non-determinism [52]. In addition, system wide page coloring introduces constraints on memory management which may interfere with the application's memory requirements [70].

Implementing page placement policies is non-trivial: The complexity of the FreeBSD kernel is increased significantly [31], Solaris allows applications to chose from multiple algorithms [61], and there have been several failed attempts to implement page placement algorithms in Linux. Other systems like COLORIS [69] replace Linux' page allocator entirely in order to support page coloring.

In contrast, Cichlid allows an application to explicitly request physical memory of a certain color and map according to its needs. For instance, a streaming database join operator can restrict the large relation (which is streamed from disk) to a small portion of the cache as most accesses would result in a cache miss anyway and keep the smaller relation completely in cache.

Table 8 shows the results of parallel execution of two instances of the HPC Challenge suite RandomAccess benchmark on cores that share the same last-level cache. In the first column we show the performance of each instance running in isolation. We see a significant drop in GUP/s for the instance with the smaller working set when both instances run in parallel. By applying cache partitioning we can keep the performance impact on the smaller instance to a minimum while improving the performance of the larger instance even compared to the case where the larger instance runs in isolation.

The reason behind this unexpected performance improvement is that the working set (the table) of the larger instance is restricted to a small fraction of the cache which reduces conflict misses between the working set and other data structures such as process state etc.

## 4.8 Discussion

With this evaluation, we have shown that the flexibility of Cichlid's memory system allows applications to opti-

| Process | Isolation | Parallel | | Parallel Colors | |
|---|---|---|---|---|---|
| 16M Table | 0.0926 | 0.0834 | 90.0% | 0.0921 | 99.5% |
| 64M Table | 0.0570 | 0.0561 | 98.4% | 0.0631 | 110.7% |

Table 8: Parallel execution of GUPS on Cichlid with and without cache coloring. Values in GUP/s.

mize their physical resources for a particular workload independent of a system-wide policy without sacrificing performance.

Cichlid's strength lies in its flexibility. By stripping back the policies baked into traditional VM systems over the years (many motivated by RAM as a scarce resource) and exposing hardware resources securely to programs, it performs as well as or better than Linux for most benchmarks, while enabling performance optimizations not previously possible in a clean manner.

## 5   Related work

Prior to Barrelfish and seL4, the idea of moving memory management into the application rather than a kernel or external paging server had been around for some time. Engler et al. in 1995 [35] outlined much of the motivation for moving memory management into the application rather than the kernel or external paging server, and described AVM, an implementation for the Exokernel [36] based on a software-loaded TLB, presenting a small performance evaluation on microbenchmarks. AVM referred to physical memory explicitly by address, and "secure bindings" conferred authorization to map it. Since then, software-loaded TLBs have fallen out of favor due to hardware performance trends. Cichlid targets modern hardware page tables, and uses capabilities to both name and authorize physical memory access.

The V++ Cache Kernel [18] implemented user-level management of physical memory through page-frame caches [46] allowing applications to monitor and control the physical frames they have, with a focus on better page-replacement policies. A virtual address space is a segment which is composed of regions from other segments called bound regions. A segment manager, associated with each segment, is responsible for keeping track of the segment to page mappings and hence handling page faults. Pages are migrated between segments to handle faults. Segment managers can run separately from the faulting application. It is critical to avoid double faults in the segment manager. Initialization is handled by the kernel which creates a well-known segment.

Other systems have also reflected page faults to user space. Microkernels like L4 [57], Mach [64], Chorus [1], and Spring [51] allow server processes to implement cus-

tom page management policies. In contrast, the soft-realtime requirements of continuous media motivated Nemesis [44] redirecting faults to the application itself, to ensures resource accountability. As with AVM, the target hardware is a uniprocessor with a software-loaded TLB. A similar upcall mechanism for reflecting page faults was used in K42 [55].

In contrast, extensible kernels like SPIN [10] and VINO [34] allow downloading of safe policy extensions into the kernel for performance. For example, SPIN's kernel interface to memory has some similarity with Cichlid's user-space API: `PhysAddr` allowed allocation, deallocation, and reclamation of physical memory, `VirtAddr` managed a virtual address space, and `Translation` allowed the installation of mappings between the two, as well as event handlers to be installed for faults. In comparison, Cichlid allows applications to define policies completely in user-space, whereas SPIN has to rely on compiler support to make sure the extensions are safe for use in kernel-space.

## 6   Conclusion

Cichlid inverts the classical VM model and securely exposes physical memory and MMU hardware to applications without recourse to virtualization hardware. It enables a variety of optimizations based on the memory system which are either impossible to express in Unix-like systems, or can only be cast as "hints" to a fixed kernel policy. Although MMU hardware has evolved to support a Unix-oriented view of virtual memory, Cichlid outperforms the Linux VM in many cases, and equals it in others.

Cichlid explores a very different style of OS service provision. Demand paging often badly impacts modern applications that rely on fast memory; the virtual address space can be an abstraction barrier which degrades performance. In Cichlid, in contrast, an application *knows* when it has insufficient physical memory and must explicitly deal with it. Given current trends in both applications and hardware, we feel this "road less travelled" in OS design is worthy of further attention. Exposing hardware securely to applications, libraries, and language runtimes may be the only practical way to avoid the increasing complexity of memory interfaces based purely on virtual addressing.

## References

[1] ABROSSIMOV, E., ROZIER, M., AND SHAPIRO, M. Generic Virtual Memory Management for Operating System Kernels. In *Proceedings of the Twelfth*

*ACM Symposium on Operating Systems Principles*
(1989), SOSP '89, ACM, pp. 123–136.

[2] AHN, J., JIN, S., AND HUH, J. Revisiting Hardware-assisted Page Walks for Virtualized Systems. In *Proceedings of the 39th Annual International Symposium on Computer Architecture* (Washington, DC, USA, 2012), ISCA '12, IEEE Computer Society, pp. 476–487.

[3] APPEL, A. W., AND LI, K. Virtual Memory Primitives for User Programs. In *Proceedings of the Fourth International Conference on Architectural Support for Programming Languages and Operating Systems* (New York, NY, USA, 1991), ASPLOS IV, ACM, pp. 96–107.

[4] ARM LTD. *Cortex-A9 Technical Reference Manual.* Revision r4p1.

[5] AZIZ, K. Improving the Performance of Transparent Huge Pages in Linux. https://blogs.oracle.com/linuxkernel/entry/performance_impact_of_transparent_huge, Aug 2014.

[6] BARR, T. W., COX, A. L., AND RIXNER, S. Translation Caching: Skip, Don't Walk (the Page Table). In *Proceedings of the 37th Annual International Symposium on Computer Architecture* (New York, NY, USA, 2010), ISCA '10, ACM, pp. 48–59.

[7] BASU, A., GANDHI, J., CHANG, J., HILL, M. D., AND SWIFT, M. M. Efficient Virtual Memory for Big Memory Servers. In *Proceedings of the 40th Annual International Symposium on Computer Architecture* (New York, NY, USA, 2013), ISCA '13, ACM, pp. 237–248.

[8] BAUMANN, A., BARHAM, P., DAGAND, P.-E., HARRIS, T., ISAACS, R., PETER, S., ROSCOE, T., SCHÜPBACH, A., AND SINGHANIA, A. The Multikernel: a new OS architecture for scalable multicore systems. In *Proceedings of the 22nd ACM Symposium on Operating Systems Principles* (2009), pp. 29–44.

[9] BELAY, A., BITTAU, A., MASHTIZADEH, A., TEREI, D., MAZIÈRES, D., AND KOZYRAKIS, C. Dune: safe user-level access to privileged CPU features. In *Proceedings of the 10th USENIX conference on Operating Systems Design and Implementation (OSDI)* (Hollywood, CA, USA, 2012).

[10] BERSHAD, B. N., SAVAGE, S., PARDYAK, P., SIRER, E. G., FIUCZYNSKI, M. E., BECKER, D., CHAMBERS, C., AND EGGERS, S. Extensibility Safety and Performance in the SPIN Operating System. In *Proceedings of the Fifteenth ACM Symposium on Operating Systems Principles* (New York, NY, USA, 1995), SOSP '95, ACM, pp. 267–283.

[11] BHARGAVA, R., SEREBRIN, B., SPADINI, F., AND MANNE, S. Accelerating Two-dimensional Page Walks for Virtualized Systems. In *Proceedings of the 13th International Conference on Architectural Support for Programming Languages and Operating Systems* (2008), ASPLOS XIII, pp. 26–35.

[12] BHATTACHARJEE, A. Large-reach Memory Management Unit Caches. In *Proceedings of the 46th Annual IEEE/ACM International Symposium on Microarchitecture* (New York, NY, USA, 2013), MICRO-46, ACM, pp. 383–394.

[13] BOEHM, H.-J. Conservative GC algorithmic overview. http://www.hboehm.info/gc/gcdescr.html.

[14] BOEHM, H.-J. Gcbench. http://hboehm.info/gc/gc_bench/.

[15] BOEHM, H.-J., DEMERS, A. J., AND SHENKER, S. Mostly Parallel Garbage Collection. In *Proceedings of the ACM SIGPLAN 1991 Conference on Programming Language Design and Implementation* (1991), PLDI '91, pp. 157–164.

[16] BOYD-WICKIZER, S., CHEN, H., CHEN, R., MAO, Y., KAASHOEK, F., MORRIS, R., PESTEREV, A., STEIN, L., WU, M., DAI, Y., ZHANG, Y., AND ZHANG, Z. Corey: An Operating System for Many Cores. In *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation* (Berkeley, CA, USA, 2008), OSDI'08, USENIX Association, pp. 43–57.

[17] CASEY, M. Performance Issues with Transparent Huge Pages (THP). https://blogs.oracle.com/linux/entry/performance_issues_with_transparent_huge, Sep 2013.

[18] CHERITON, D. R., AND DUDA, K. J. A Caching Model of Operating System Kernel Functionality. In *Proceedings of the 1st USENIX Conference on Operating Systems Design and Implementation* (Monterey, California, 1994), OSDI '94, USENIX Association.

[19] CORBET, J. AutoNUMA: the other approach to NUMA scheduling. http://lwn.net/Articles/488709/, Mar 2012.

[20] CORBET, J. NUMA in a hurry. http://lwn.net/Articles/524977/, Nov 2012.

[21] CORBET, J. Toward better NUMA scheduling. http://lwn.net/Articles/486858/, Mar 2012.

[22] CORBET, J. NUMA scheduling progress. http://lwn.net/Articles/568870/, Oct 2013.

[23] CORBET, J. User-space page fault handling. http://lwn.net/Articles/550555/, May 2013.

[24] CORBET, J. 2014 LSFMM summit: Huge page issues. http://lwn.net/Articles/592011/, Mar 2014.

[25] CORBET, J. NUMA placement problems. http://lwn.net/Articles/591995/, Mar 2014.

[26] CORBET, J. Page faults in user space: MADV_USERFAULT, remap_anon_range(), and userfaultfd(). http://lwn.net/Articles/615086/, Oct 2014.

[27] CORBET, J. Transparent huge pages in 2.6.38. http://lwn.net/Articles/423584/, Jan 2014.

[28] DAGAND, P.-E., BAUMANN, A., AND ROSCOE, T. Fileto-o-Fish: practical and dependable domain-specific languages for OS development. In *5th Workshop on Programming Languages and Operating Systems (PLOS)* (Oct 2009).

[29] DASHTI, M., FEDOROVA, A., FUNSTON, J., GAUD, F., LACHAIZE, R., LEPERS, B., QUEMA, V., AND ROTH, M. Traffic Management: A Holistic Approach to Memory Placement on NUMA Systems. In *Proceedings of the Eighteenth International Conference on Architectural Support for Programming Languages and Operating Systems* (Houston, Texas, USA, 2013), ASPLOS '13, ACM, pp. 381–394.

[30] DERRIN, P., ELKADUWE, D., AND ELPHINSTONE, K. *seL4 Reference Manual*. NICTA, 2006. http://www.ertos.nicta.com.au/research/sel4/sel4-refman.pdf.

[31] DILLON, M. Design elements of the FreeBSD VM system - Page Coloring. Online, https://www.freebsd.org/doc/en/articles/vm-design/page-coloring-optimizations.html, Nov 2013. Accessed 2015-08-26.

[32] ELKADUWE, D., DERRIN, P., AND ELPHINSTONE, K. A memory allocation model for an embedded microkernel. In *Proceedings of the 1st International Workshop on Microkernels for Embedded Systems (MIKES)* (2007), pp. 28–34.

[33] ELKADUWE, D., DERRIN, P., AND ELPHINSTONE, K. Kernel Design for Isolation and Assurance of Physical Memory. In *Proceedings of the 1st Workshop on Isolation and Integration in Embedded Systems* (New York, NY, USA, 2008), IIES '08, ACM, pp. 35–40.

[34] ENDO, Y., SELTZER, M., GWERTZMAN, J., SMALL, C., SMITH, K. A., AND TANG, D. VINO: The 1994 Fall Harvest. Technical Report TR-34-94, Center for Research in Computing Technology, Harvard University, December 1994.

[35] ENGLER, D. R., GUPTA, S. K., AND KAASHOEK, M. F. AVM: Application-level Virtual Memory. In *Proceedings of the Fifth Workshop on Hot Topics in Operating Systems (HotOS-V)* (1995), HOTOS '95, IEEE Computer Society, pp. 72–.

[36] ENGLER, D. R., KAASHOEK, M. F., AND O'TOOLE, JR., J. Exokernel: An Operating System Architecture for Application-level Resource Management. In *Proceedings of the 15th ACM Symposium on Operating Systems Principles* (1995), pp. 251–266.

[37] EVANS, J. Issue #243: Improve interaction with transparent huge pages. https://github.com/jemalloc/jemalloc/issues/243, Jul 2015.

[38] GAUD, F., LEPERS, B., DECOUCHANT, J., FUNSTON, J., FEDOROVA, A., AND QUÉMA, V. Large Pages May Be Harmful on NUMA Systems. In *Proceedings of the 2014 USENIX Conference on USENIX Annual Technical Conference* (Philadelphia, PA, 2014), USENIX ATC'14, USENIX Association, pp. 231–242.

[39] GICEVA, J., ALONSO, G., ROSCOE, T., AND HARRIS, T. Deployment of query plans on multicores. *Proc. VLDB Endow. 8*, 3 (Nov 2014), 233–244.

[40] GORMAN, M. Huge pages. http://lwn.net/Articles/374424/, Feb 2010.

[41] GORMAN, M. Huge pages part 2: Interfaces. https://lwn.net/Articles/375096/, Feb 2010.

[42] GORMAN, M., AND HEALY, P. Performance Characteristics of Explicit Superpage Support. In *Proceedings of the 2010 International Conference on Computer Architecture* (Berlin, Heidelberg, 2012), ISCA'10, Springer-Verlag, pp. 293–310.

[43] HAIBLE, B., AND BONZINI, P. GNU libsigsegv - Handling page faults in user mode. http://libsigsegv.sourceforge.net/.

[44] HAND, S. M. Self-paging in the Nemesis Operating System. In *Proceedings of the Third Symposium on Operating Systems Design and Implementation* (New Orleans, Louisiana, USA, 1999), OSDI '99, USENIX Association, pp. 73–86.

[45] HANSEN, D. TLB flushing on x86. https://www.kernel.org/doc/Documentation/x86/tlb.txt.

15

[46] HARTY, K., AND CHERITON, D. R. Application-controlled Physical Memory Using External Page-cache Management. In *Proceedings of the Fifth International Conference on Architectural Support for Programming Languages and Operating Systems* (New York, NY, USA, 1992), ASPLOS V, ACM, pp. 187–197.

[47] HONG, S., CHAFI, H., SEDLAR, E., AND OLUKOTUN, K. Green-Marl: A DSL for Easy and Efficient Graph Analysis. In *Proceedings of the Seventeenth International Conference on Architectural Support for Programming Languages and Operating Systems* (New York, NY, USA, 2012), ASPLOS XVII, ACM, pp. 349–362.

[48] HP LABS. The Machine. http://www.hpl.hp.com/research/systems-research/themachine/, January 2015.

[49] INTEL CORPORATION. *Intel 64 and IA-32 Architectures Optimization Reference Manual*, September 2014. Online. Accessed 2015-03-12. http://www.intel.com/content/www/us/en/architecture-and-technology/64-ia-32-architectures-optimization-manual.html?wapkw=order+number+248966-025.

[50] KAESTLE, S., ACHERMANN, R., ROSCOE, T., AND HARRIS, T. Shoal: Smart Allocation and Replication of Memory for Parallel Programs. In *Proceedings of the 2015 USENIX Annual Technical Conference* (Santa Clara, CA, 2015), USENIX ATC '15, pp. 263–276.

[51] KHALIDI, Y. A., AND NELSON, M. N. The Spring Virtual Memory System. Technical Report SMLI TR-93-9, Sun Microsystems Laboratories Inc., February 1993.

[52] KIM, J., KIM, J., AHN, D., AND EOM, Y. I. Page coloring synchronization for improving cache performance in virtualization environment. In *Computational Science and Its Applications-ICCSA 2011.* Springer, 2011, pp. 495–505.

[53] KLEIN, G., ELPHINSTONE, K., HEISER, G., ANDRONICK, J., COCK, D., DERRIN, P., ELKADUWE, D., ENGELHARDT, K., KOLANSKI, R., NORRISH, M., SEWELL, T., TUCH, H., AND WINWOOD, S. seL4: Formal verification of an OS kernel. In *Proceedings of the 22nd ACM Symposium on Operating Systems Principles* (2009).

[54] KOESTER, D., AND LUCAS, B. HPC Challenge - Random Access. Online. http://icl.cs.utk.edu/projectsfiles/hpcc/RandomAccess/. Accessed 2015-03-09.

[55] KRIEGER, O., AUSLANDER, M., ROSENBURG, B., WISNIEWSKI, R. W., XENIDIS, J., DA SILVA, D., OSTROWSKI, M., APPAVOO, J., BUTRICO, M., MERGEN, M., WATERLAND, A., AND UHLIG, V. K42: Building a Complete Operating System. In *Proceedings of the 1st EuroSys Conference* (2006), pp. 133–145.

[56] LEIS, V., BONCZ, P., KEMPER, A., AND NEUMANN, T. Morsel-driven parallelism: A numa-aware query evaluation framework for the many-core age. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data* (New York, NY, USA, 2014), SIGMOD '14, ACM, pp. 743–754.

[57] LIEDTKE, J., UHLIG, V., ELPHINSTONE, K., JAEGER, T., AND PARK, Y. How to Schedule Unlimited Memory Pinning of Untrusted Processes or Provisional Ideas About Service-Neutrality. In *Proceedings of the The Seventh Workshop on Hot Topics in Operating Systems* (Washington, DC, USA, 1999), HOTOS '99, IEEE Computer Society, pp. 153–.

[58] LINUX KERNEL PROJECT. Hugetlbpage support in the Linux kernel. https://www.kernel.org/doc/Documentation/vm/hugetlbpage.txt.

[59] LINUX KERNEL PROJECT. Transparent Hugepage Support. https://www.kernel.org/doc/Documentation/vm/transhuge.txt.

[60] NAVARRO, J., IYER, S., DRUSCHEL, P., AND COX, A. Practical, transparent operating system support for superpages. *SIGOPS Oper. Syst. Rev. 36*, SI (Dec 2002), 89–104.

[61] ORACLE CORPORATION. Online. http://docs.oracle.com/cd/E19683-01/806-7009/chapter2-95/index.html, 2010. Accessed 2015-08-15.

[62] PAOLO FARABOSCHI, KIMBERLY KEETON, T. M., AND MILOJICIC, D. Beyond processor-centric operating systems. In *Proceedings of the 2015 International Workshop on Hot Topics in Operating Systems (HotOS XV)* (Karthause Ittingen, Warth-Weiningen, Switzerland, May 2015).

[63] PETER, S., LI, J., ZHANG, I., PORTS, D. R. K., WOOS, D., KRISHNAMURTHY, A., ANDERSON, T., AND ROSCOE, T. Arrakis: The Operating System is the Control Plane. In *11th Symposium on Operating Systems Design and Implementation (OSDI'14)* (Broomfield, Colorado, USA, October 2014).

[64] RASHID, R., TEVANIAN, A., J., YOUNG, M., GOLUB, D., BARON, R., BLACK, D., BOLOSKY, W., AND CHEW, J. Machine-Iindependent Virtual Memory Management for Paged Uniprocessor and Multiprocessor

Architectures. *Computers, IEEE Transactions on 37,*
8 (Aug 1988), 896–908.

[65] SANFILIPPO, S. Redis latency problems troubleshoot-
ing. http://redis.io/topics/latency.

[66] Soft-Dirty PTEs. https://www.kernel.org/doc/
Documentation/vm/soft-dirty.txt.

[67] SOMA, Y., GEROFI, B., AND ISHIKAWA, Y. Revisiting
Virtual Memory for High Performance Computing
on Manycore Architectures: A Hybrid Segmenta-
tion Kernel Approach. In *Proceedings of the 4th
International Workshop on Runtime and Operating
Systems for Supercomputers* (New York, NY, USA,
2014), ROSS '14, ACM, pp. 3:1–3:8.

[68] THE UNIVERSITY OF TENNESSEE. HPC Challenge
Benchmark. Online. http://icl.cs.utk.edu/hpcc/
software/view.html?id=178. Accessed 2015-03-09.

[69] YE, Y., WEST, R., CHENG, Z., AND LI, Y. COLORIS:
A Dynamic Cache Partitioning System Using Page
Coloring. In *Proceedings of the 23rd International
Conference on Parallel Architectures and Compi-
lation* (Edmonton, AB, Canada, 2014), PACT '14,
pp. 381–392.

[70] ZHANG, X., DWARKADAS, S., AND SHEN, K. Towards
Practical Page Coloring-based Multicore Cache
Management. In *Proceedings of the 4th ACM Euro-
pean Conference on Computer Systems* (Nuremberg,
Germany, 2009), EuroSys '09, pp. 89–102.