

Twizzler: A *Data-centric* OS for Non-volatile Memory

DANIEL BITTMAN and PETER ALVARO, UC Santa Cruz, USA

PANKAJ MEHRA, IEEE Member, USA

DARRELL D. E. LONG, UC Santa Cruz, USA

ETHAN L. MILLER, UC Santa Cruz, USA and Pure Storage, USA

Byte-addressable, non-volatile memory (NVM) presents an opportunity to rethink the entire system stack. We present Twizzler, an operating system redesign for this near-future. Twizzler removes the kernel from the I/O path, provides programs with memory-style access to persistent data using small (64 bit), object-relative cross-object pointers, and enables simple and efficient long-term sharing of data both between applications and between runs of an application. Twizzler provides a clean-slate programming model for persistent data, realizing the vision of UNIX in a world of persistent RAM.

We show that Twizzler is simpler, more extensible, and more secure than existing I/O models and implementations by building software for Twizzler and evaluating it on NVM DIMMs. Most persistent pointer operations in Twizzler impose less than 0.5 ns added latency. Twizzler operations are up to 13× faster than UNIX, and SQLite queries are up to 4.2× faster than on PMDK. YCSB workloads ran 1.1–2.9× faster on Twizzler than on native and NVM-optimized SQLite backends.

CCS Concepts: • **Software and its engineering** → **Operating systems**; • **Hardware** → **Memory and dense storage**; • **Information systems** → **Storage class memory**;

Additional Key Words and Phrases: Persistent memory, non-volatile memory, NVM, PMEM, single-level store, global address space, memory hierarchy

ACM Reference format:

Daniel Bittman, Peter Alvaro, Pankaj Mehra, Darrell D. E. Long, and Ethan L. Miller. 2021. Twizzler: A *Data-centric* OS for Non-volatile Memory. *ACM Trans. Storage* 17, 2, Article 11 (June 2021), 31 pages.

<https://doi.org/10.1145/3454129>

1 INTRODUCTION

Byte-addressable **non-volatile memory (NVM)** on the memory bus with DRAM-like latency [24, 40] will fundamentally shift the way that we program computers. The two-tier memory hierarchy split between high-latency persistent storage and low latency volatile memory may evolve into a single level of large, low latency, and directly addressable persistent memory. Mere incremental change will leave dramatic improvements in programmability, performance, and simplicity on the

This work was supported in part by the National Science Foundation (grants IIP-1266400, IIP-1841545), a grant from Intel Corporation, and the industrial members of the UCSC Center for Research in Storage Systems.

Authors' addresses: D. Bittman and P. Alvaro, UC Santa Cruz, CSE, 1156 High Street, Santa Cruz, CA; emails: {dbittman, palvaro}@ucsc.edu; P. Mehra, IEEE Member; emails: pankaj.mehra@ieee.org; D. D. E. Long, UC Santa Cruz, BSOE, 1156 High Street, Santa Cruz, CA; email: darrell@ucsc.edu; E. L. Miller, UC Santa Cruz, CSE, 1156 High Street, Santa Cruz, CA, Pure Storage, CA; email: elm@ucsc.edu.

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

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.

1553-3077/2021/06-ART11 \$15.00

<https://doi.org/10.1145/3454129>

table. It is essential that operating systems and system software evolve to make the best use of this new technology.

These opportunities motivate us to revisit how programs operate on persistent data. The separation of volatile memory and high-latency persistent storage is a core OS design principle that requires the OS to manage ephemeral copies of data and interpose itself on persistence operations, a penalty that will consume an increasing fraction of time as NVM performance increases [65]. The direct-access nature of NVM invites the use of load and store instructions to directly access persistent data, simplifying applications by enabling persistent data manipulation without the need to transform it between in-memory and on-storage data formats. Thus, the model that best exploits the low latency nature of NVM is one in which persistent data is maintained as in-memory data structures and not serialized or explicitly loaded or unloaded. To avoid serialization, this model must support *persistent pointers* that are valid in *any* execution context, not just the one in which they were created.

Trying to mold NVM into existing models will not enable its fullest potential, just as SSDs did not reach their full potential until they transcended the disk paradigm. To explore a “clean-slate” approach, we are building Twizzler, an OS designed to take full advantage of this new technology by rethinking the abstractions OSES provide in the context of NVM. Twizzler divides NVM into *objects* within a global object space, and pointers are interpreted in the context of the object in which they reside. This decouples pointers from the address space of an individual thread, providing a data-centric programming model rather than a process-centric one. The result is a vastly simpler environment in which the OS’s primary function is to support manipulating, sharing, and protecting persistent data using few kernel interpositions.

We designed and implemented a simple, standalone kernel that supports a userspace for NVM-based applications, with compatibility layers for legacy programs. We wrote a set of libraries and portability layers that provide a rich environment for applications to access persistent data that takes into account both semantics (persistent pointers) and safety (building crash-consistent data structures). We then performed a case-study by writing software for Twizzler, taking into account the new flexibility and power gained by our model and evaluating our software for complexity and performance. We ported SQLite to Twizzler, showing how our approach can provide significant performance gains on existing applications as well.

In a world where in-memory data can last forever, the context required to manipulate that data is best coupled with the *data* rather than ephemeral constructs like the process. This key insight manifests itself in the three primary contributions of this article:

- We discuss (Section 2) our vision for a data-centric OS and the requirements that it must meet to provide low latency memory-style access to NVM with efficient data sharing and simplified programming models.
- We present Twizzler (Section 3) and describe its mechanisms to meet those requirements, including decoupling traditionally linked concerns, reducing kernel involvement in address space management, and providing a rich model for constructing in-memory persistent data structures that can be easily shared between programs and machines.
- We evaluate (Section 5) the ease-of-use, security advantages, and programmability offered by our environment, for both new and existing, ported software (SQLite), along with performance improvements (Section 6) on NVM DIMMs.

2 THE DATA-CENTRIC OS

Operating systems provide abstractions for data access that reflect the hardware for which they were designed. Current I/O interfaces and abstractions reflect the structure of mutually exclusive

volatile and persistent domains, the hallmarks of which are heavy kernel involvement for persisting data, a need for data serialization, and complexity in data sharing requiring the overhead of pipes or the management cost of shared virtual memory. However, the introduction of low latency and directly attached NVM into the memory hierarchy requires that we rethink key assumptions such as the use of virtual addresses, the kernel’s involvement in persistent I/O, and the way that programs operate on and share persistent data [31].

The first key characteristic of NVM is low latency: only 1.5–8× the latency of DRAM in most cases [40]. Thus, the cost of a system call to access NVM dominates the latency of the access itself. The second key characteristic is that the processor can directly access persistent storage using load and store instructions. Direct, low latency access to NVM means that explicit serialization is a poor fit—it adds complexity, as programmers must maintain different data formats and the transformations between them, and the overhead is intolerable due to NVM’s low latency. Hence, we should design the semantics of the programming model around *in-memory* persistent data structures, giving programs direct access to them without explicit persistence calls or serialization methods.

These characteristics imply two basic requirements for OSes to most effectively use NVM:

- (R1) **Remove the kernel from the persistence path.** This addresses both characteristics. System calls to persist data are costly; we must provide lightweight, direct, memory-style access for programs to operate on persistent data.
- (R2) **Design for pointers that last forever.** Long-lived data structures can directly reference persistent data, so pointers must have the same lifetime as the data they point to. Virtual memory mappings are, by contrast, ephemeral and so cannot effectively name persistent data. Persistent data is, by definition, accessed by multiple actors, both simultaneously and over time, and thus must be stored in a form that is conducive to sharing without needing the ephemeral context associated with a particular actor.

We call an OS that meets both requirements R1 and R2 *data-centric*, as opposed to current OSes, which are *process-centric*. Operations on persistent, in-memory data structures are the primary functions of a data-centric OS, which tries to avoid interposing on such operations, preferring instead to intervene only when necessary to ensure properties such as security and isolation. To meet both of these requirements, a data-centric OS must provide effective abstractions for identifying data independent of data location, constructing persistent data relationships that do not depend on ephemeral context, and facilitating sharing and protection of persistent data.

2.1 A Data-centric Approach

We cannot store virtual addresses in persistent data, so we need a new way to name a word of persistent memory: a *persistent pointer*. The persistent pointer encodes a persistent identification of data (Section 3.3) instead of an ephemeral address, allowing any thread to access the desired word of memory regardless of address space. This approach dramatically improves programmability, as programmers need not worry about the complexity of referring to persistent data with ephemeral constructs, improving data sharing between programs and across runs of a program. Twizzler still makes use of virtual memory *hardware* to provide isolation and translation, but persistent data structures should not be written in terms of virtual addresses.

The Death of the Process. Processes as a first class OS abstraction are, like virtual addresses, unnecessary; a traditional process couples threads of control to a virtual address space, a security role, and kernel state. However, with the kernel removed from persistent data access, much of that kernel state (e.g., file descriptors) is unnecessary, leading to a decoupling of mechanisms: Nothing fundamentally connects a virtual address space (a piece of ephemeral context used to access data)

and a security context (*what* data threads may access). Instead, a data-centric OS can keep the good parts of a process but *separate* virtual address translation and security roles, allowing threads to select one of each as needed.

The process abstraction is just one example. Persistent data access plays a key role in OS abstraction design, and we need to avoid complexity arising from combining old and new interfaces. Hence, we need to consider the wide-reaching effects of changing the persistence model on *all* aspects of the system, not just I/O interfaces. NVM gives us an opportunity to design an OS around the requirements of the target programming model instead of trying to mold support libraries around existing interfaces. While it is important that we provide support for legacy applications, it is these applications that should be relegated to support libraries; new applications built for the programming model should get first-class OS support.

Targeting These Constraints with Twizzler. The consequences of meeting the requirements of these hardware trends define a bounded design space for data-centric OSes. We have chosen a point in that space and built Twizzler, our approach to providing applications with efficient and effective access to NVM. In the following section, we will discuss how our four primary abstractions—a low-level persistent object model, a persistent pointer design, an address space mechanism called *views*, and a security context mechanism—achieve these goals of removing the kernel from the persistent data access path.

2.2 Existing Interfaces

Current OS techniques do not meet requirements R1 and R2 as we set out above—file read and write interfaces, designed for sequential media and later expanded for block-based media, require significant kernel involvement and serialization, violating both requirements. While support for these interfaces can be useful for legacy applications, as we will demonstrate, providing the programmer with abstractions designed *for* NVM both reduces complexity and improves performance.

The `mmap` system call attempts to hide storage behind a memory interface through hidden data copies. But, with NVM, these copies are wasteful, and `mmap` still has significant kernel involvement and the need for explicit `msync` calls. “**Direct Access**” (**DAX**) tries to retrofit `mmap` for NVM by removing the redundant copy, but this *still* fails to address requirement R2! Operating on persistent data through `mmap` requires the programmer to use either fixed virtual addresses, which presents an infeasible coordination problem as we scale across machines, or virtual addresses directly, which are ephemeral and require the context of the process that created them.

Attempting to shoehorn NVM programming atop POSIX interfaces (including `mmap`) results in complexity that arises from combining multiple partial solutions. Given some feature desired by an application, the NVM framework can provide an integrated solution that meshes well with the existing support for persistent data structure manipulation and access, or it can fall back to POSIX, resulting in the programmer needing to understand two different “feature namespaces” and their interactions. An example of this is naming, where a programmer may need to turn to the filesystem to manage names in a completely orthogonal way to how the NVM frameworks handles data references. For example, PMDK, an NVM programming library, relies on a filesystem for naming and initial access to persistent memory objects, resulting in different kinds of references, feature sets from filesystems being applied (like security) while others are not (data access), and the complexity of understanding how the PMDK abstractions interact with the POSIX ones. Instead, our model prefers to build legacy support atop new abstractions (Section 4) and avoid falling back to legacy models for persistent data access. We will discuss another example, security, in our case study (Section 5).

Additionally, models that layer NVM programming atop existing interfaces often fail to facilitate effective persistent data sharing and protection. PMDK, for example, makes design choices that

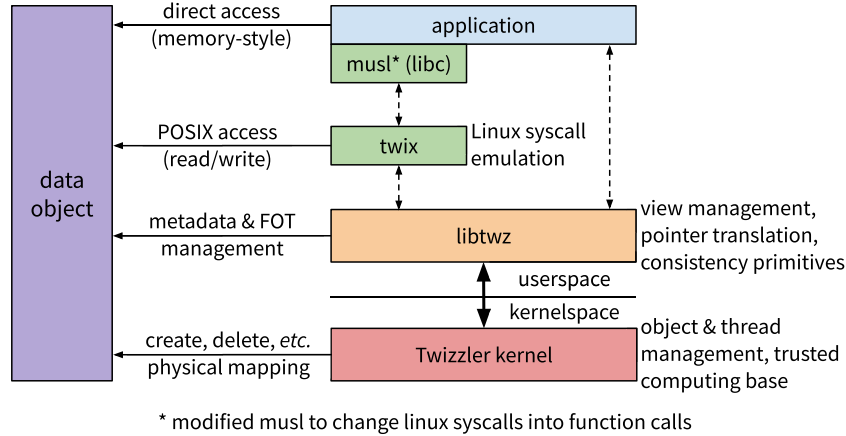


Fig. 1. Twizzler system overview. Applications link to `musl` (a C library), `twix` (our Linux syscall emulation library), and `libtwz` (our standard library). Through `musl`, they may act on persistent data with POSIX interfaces, though we expect Twizzler applications to operate directly on persistent data with memory-style semantics.

limit scalability, since its data objects are not self-contained and do not have a large enough ID space, resulting in the need to coordinate object IDs across machines [11]. For the same reason, although single-address space OSes [13] somewhat address requirement R1, they do not consider both requirements at once, nor do they provide an effective and scalable solution to long-term data references due to that same coordination complexity [10].

3 THE DESIGN OF TWIZZLER

Twizzler is a stand-alone kernel and userspace runtime that provides execution support for programs. It provides, as first-class abstractions, a notion of threads, address spaces, persistent objects, and security contexts. A program typically executes as a number of threads in a single address space (providing backwards compatibility with existing programming models), into which persistent objects are mapped on-demand. Instead of providing a process abstraction, Twizzler provides *views* (Section 3.2) of the object space, which formalizes the notion of ephemeral context within our model by allowing programs to map objects for access, and *security contexts* (Section 3.4), which define a thread's access rights to objects in the system. Twizzler provides persistent pointers (Section 3.3) for programs, as well as primitives to ensure crash-consistency (Section 4.2). The thread abstraction is similar to modern OSes: The kernel provides scheduling, synchronization, and management primitives. Figure 1 shows an overview of the system organization and how different parts of the system operate on data objects.

Twizzler's kernel acts much like an Exokernel [29, 43], providing sufficient services for a userspace library OS, called `libtwz`, to provide an execution environment for applications. The primary job of `libtwz` is to manage mappings of persistent objects into the address space (Section 3.2) and deal with persistent pointers (Section 3.3). Twizzler also exposes a standard library that provides higher-level interfaces beyond raw access to memory. For example, software that better fits message-passing semantics can use library routines that implement message-passing atop shared memory. Twizzler's standard library provides additional higher-level interfaces, including streams, logging, event notification, and many others. Applications use these to easily build composable tools and pipelines for operating on in-memory data structures without the performance loss and complexity of explicit I/O.

We provide POSIX support with `twix`, a library that emulates Linux syscalls. We modified `musl` [1], a C library that all programs link to, replacing invocations of the `syscall` instruction with calls into `twix`, which internally tracks UNIX state-like file descriptors. This is handled entirely in userspace; calls to `read` and `write` often reduce to calls to `memcpy`.

3.1 Object Management

Twizzler organizes data into *objects*, which may be persistent. Each object is identified by a unique 128 bit object ID (though larger IDs would be possible). Objects provide contiguous regions of memory that organize semantically related data with similar lifetime and permissions. Applications access objects via mapping services (discussed in the next section) by mapping each object into a contiguous range in the address space, though the address space itself may be densely or sparsely mapped. Objects can be anywhere from 4 KiB (the size of a page) to 1 GiB; the upper bound on object size is a prototype implementation choice and not fundamental to the design.

Twizzler uses objects as the unit of access control, building off a read/write/execute permissions model that mirrors that of memory management units in modern processors. This is a direct consequence of avoiding the kernel for persistent data access—it can set policy by programming the MMU, but must leave enforcement up to the hardware, which, in turn, defines what protections are possible.

An object, from the programmer’s perspective, is flexible in its contents—for example, it could contain anywhere from a single B-tree node to the entire B-tree. Often, an object would contain the entire tree, since the entire tree is typically subject to the same access semantics by programs, and there are overheads associated with objects that can be amortized over larger spaces. Data and data structures that are too large for one object or require different access permissions can span multiple objects with references between them. We demonstrate the benefits of this flexibility in Section 5.

The kernel provides services for object management, such as creating and deleting objects. Objects are created by the `create` system call, which returns an object ID. A program may also optionally provide an existing object ID to the `create` call, stating that the new object should be a copy of the existing one, for which Twizzler uses copy-on-write. The new ID is a number that is unlikely to collide with existing IDs in the 128 bit ID space and can be assigned using a technique that supports this requirement (random, hashing, etc.). Some forms of ID assignment support a form of access control: A program can only access an object whose ID it knows. Twizzler provides object naming as well, discussed in Section 3.3.

Objects may be deleted via the `delete` system call. Like UNIX’s `unlink`, objects are reference counted, where a reference refers to a mapping in an address space. Once the reference count reaches zero, the object may be deleted. During deletion, an object may be optionally marked as “hidden,” causing new mapping requests for this object to fail.

Object Types, Persistence, and Lifetime. Applications need to be able to specify what *type* of memory an object resides in. Currently, we are operating on systems that contain both persistent NVM and volatile DRAM as main memory, and applications may want to make use of both of these memory types. Placing certain objects in DRAM, for example, can result in performance improvements (e.g., caching read-only objects) or security improvements (e.g., making temporary key material volatile). Twizzler exposes this choice to applications at object creation time, allowing them to specify the type of the object. At least two types, *volatile* and *persistent*, are supported by default. As additional types of physical memory are added to systems (e.g., different kinds of NVM with different properties, high-bandwidth memory), applications may wish to have more fine-grained control over where objects are placed, and Twizzler’s APIs allow such control. Objects can also be

moved between types of memory after creation, though this may be a time-consuming operation, as it involves copying potentially large amounts of data.

By default, objects are persistent and live in kernel-managed NVM unless they are marked as volatile. If an object is volatile, then it has a limited lifetime that is related to the power state of the machine—as soon as power is lost (or the system is rebooted) all volatile objects disappear. Note that Twizzler removes the distinction between volatile and persistent objects for how applications access data, relying on higher-level language or library support and application support for dealing with the limited lifetime of volatile objects.

The property of persistent versus volatile for objects differs from the concept of ephemeral data. The “volatile” property places a physical restriction on the lifetime of an object (the machine’s power state), while the “persistent” property indicates that the object will exist until explicitly deleted. Objects can also be long-lived or ephemeral, independent of their persistence property, since we use the term “ephemeral” to describe information, data, or state that has a finite lifetime and is expected to “go away.” While all volatile objects are ephemeral, the reverse is not true—we may place ephemeral data in a persistent object to allow for recovery after an unexpected power cycle. The “persistent” property of an object is a recorded piece of information that the kernel associates with an object, but there is no such information for ephemeral versus long-lived. Instead, we provide a mechanism for specifying a logical lifetime of objects relative to one another with a mechanism called *ties*, which we will discuss below.

Object Ties and Logical Lifetime. Applications in Twizzler also have some lifetime; an application’s job is typically to operate on some persistent data while performing some computation before eventually exiting. Such an application will likely use volatile objects to represent temporary computation state (e.g., the stack and heap, which are ephemeral). However, just assigning an object as volatile is insufficient, because there is a lifetime mismatch: The volatile object will live until the next reboot, while the application may exit before then or may even live and try to recover after a power cycle. Simply manually deleting the volatile object when the application is done is also insufficient, as it does not account for crashes where the application may be unable to clean up its state. Furthermore, applications that wish to support recovery may make use of persistent stacks and heaps, thus these objects would have to be persistent despite being ephemeral.

While we could provide a mechanism designed specifically for this “system-level” task, where the kernel maintains a set of objects to automatically cleanup when an application exits, this would require the kernel to have some understanding of what an “application” is. Furthermore, if we generalize a solution to automatic cleanup, then we can allow applications to make use of it for their own purposes. For example, in UNIX, it is common for programs to create and immediately unlink files to ensure the system frees those resources when the program exits. We would like to reproduce similar semantics here that also solves the lower-level problem above of freeing application state by assigning a lifetime to objects that is more expressive than simply “volatile” and “persistent.”

In Twizzler, object lifetime is expressed through *ties*. An object can be tied to another by invoking a system call that tells the kernel that object A is tied to object B, after which the lifetime of A is guaranteed to be at least that of B. The kernel will not fully delete object A (even if the delete system call is invoked on it) until after B is fully deleted. An object may be tied to a large (but finite) number of other objects and may also be *untied* at any time. This model of specifying object lifetime relative to others is similar to Rust [2], where reference lifetime can be named so the programmer can express lifetimes of objects relative to each other. Note that object ties are not related to persistent pointers (discussed in more detail in Section 3.3), and instead primarily provide a way to formalize automatic cleanup.

Object ties provide a convenient mechanism for applications to build large data structures across multiple objects without giving up easy cleanup if something goes wrong or if the “root” object is deleted. Twizzler also uses ties internally: When an object is created as copy-from an existing object, it uses copy-on-write semantics, and thus internally marks the source object as tied to the new object. We also tie ephemeral program state objects to threads (which are also represented by objects) such that they are automatically cleaned up when a program exits. It is our expectation that application programmers will only rarely directly use ties. Instead, we expect that ties will provide necessary features that higher-level programming language support for persistent memory can use.

Note that object ties interact with the notion of volatile and persistent objects, because volatile objects have an implicit *maximum* lifetime—that of the next machine restart or power loss. Tying volatile objects to volatile objects and persistent objects to persistent objects both act as expected. Tying a persistent object to a volatile object is also semantically simple (persistent objects already have an “assumed lifetime” that is longer than a volatile object). Tying a volatile object to a persistent object, however, may seem somewhat nonsensical. However, Twizzler does still allow this, because it has useful semantics: If an application creates a data structure with some volatile component,¹ then it may want to tie the lifetime of that volatile component to the persistent component if the data structure is to be deleted. This use case (creating a persistent object that we expect to delete) is not uncommon, particularly in applications designed to recover partial computation after a crash. Note that, in this case, the maximum lifetime of the volatile object is still in play; after a power cycle, that object will no longer be present, so tying a volatile object to a persistent object is somewhat dangerous.

3.2 Ephemeral State Management with Views

Despite Twizzler’s focus on persistent data, many components of our hardware and applications are built around ephemeral constructs. For example, threads are ephemeral “moments of computation” that act on persistent data, while the programs that they execute often expect some ephemeral private data (e.g., the data segment and the stack). While virtual addresses are the wrong abstraction for persistent data access, modern hardware provides (and often requires) the use of virtual address hardware that we can leverage for protection and isolation, adding additional ephemeral state.

Twizzler defines objects called “views,” which coalesce the state and context necessary to support ephemeral constructs like threads and application instances into Twizzler objects. A significant part of that state is ephemeral virtual address mappings; Twizzler provides access to persistent objects by mapping them into the virtual address space behind-the-scenes (via `libtwz`). The view object contains structures to define the layout of the virtual address space that the kernel reads and uses to program the MMU accordingly. Figure 2 shows how views “mesh” ephemeral threads with persistent data by providing them a context to operate in. Since view objects are normal Twizzler objects, they can be persisted, allowing us to recover application state after power cycles.

By coalescing this ephemeral state into an object, we make it possible for applications to manage it directly with minimal kernel involvement. Avoiding the kernel is natural—all data access already does this in Twizzler, so adding a separate kernel API to manage this state would add complexity—and reduces the number of system calls needed when mapping objects. Additionally, avoiding

¹Since Twizzler’s kernel is not involved in reference creation, it cannot prevent such a reference from being created. We expect language support for persistent data structures to impose restrictions on applications in this regard, and the OS should not prematurely restrict how applications use volatile and persistent objects. Access to a volatile object that no longer exists after a reboot results in a simple access fault, mitigating security concerns.

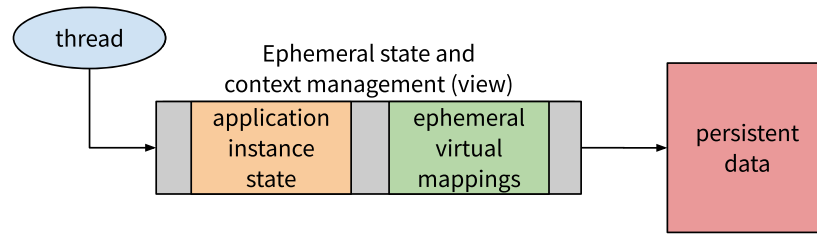


Fig. 2. View objects in Twizzler. Views manage ephemeral constructs and state, giving threads the necessary context to execute and access persistent data.

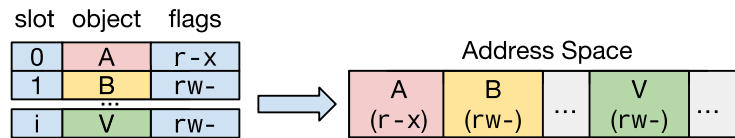


Fig. 3. Layout of a view object. The kernel consults the view object’s mapping on page-fault and maps in the requested objects at the appropriate location in the virtual address space.

the kernel necessitates an increased address space management responsibility for userspace. For example, executable loading and mapping is largely handled without the kernel.

Figure 3 shows how view objects lay out the address space of any threads running inside that particular view. View objects are manipulated by userspace and interpreted by the kernel. When applications map objects, they update the view to specify that that object should be addressable at a specific location. On a page-fault, the kernel reads the view and maps the object at the requested location. The view object is laid out like a page-table, where each entry in the table corresponds to a slot in the virtual address space. Each table entry contains an object ID and requested protection bits to further protect objects atop access control mechanisms (similar to `PROT_*` in `mmap`).

When a page-fault occurs, the fault handler tries to handle the fault by either doing copy-on-write, checking permissions, or by trying to map an object into a slot if the view object requested one. If it cannot handle the fault (due to a protection error or an empty entry in the view object), then it elevates the fault to userspace where `libtwz` handles it, possibly by killing the thread, or possibly by mapping an object if the slot is “on-demand.” This is similar to userspace paging systems [3, 35]. When the kernel maps an object into a slot, it updates the address space’s page-tables appropriately.

Applications can add objects to a view with the `view_set` function. The caller specifies a target object and a set of protections (see Section 3.3) and a slot in which to map the object. However, applications rarely invoke this function directly—instead, `libtwz` provides a higher-level API to allow applications to operate above the level of manually mapping objects. The standard library also provides access to other utility functions for views (such as querying state, creating new views, and copying views). These functions, by default, operate on a thread’s current view, but they may also optionally operate on any other view object, which allows Twizzler to implement operations with semantics similar to `fork` and `execve`.

When threads add entries to a view object they need not inform the kernel—when a fault occurs, the kernel will read the entry as needed. However, when *changing* or *deleting* an entry, threads must inform the kernel so it can update existing page table entries. We provide two system calls for views. The `become` system call allows a thread to change to a new view, which might be used to execute a new program or jump across programs to, for example, accomplish a protected task.

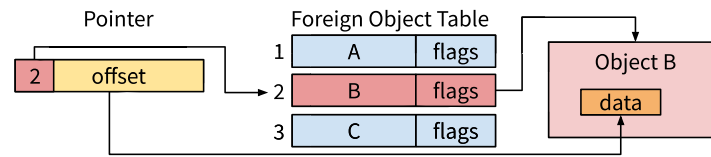


Fig. 4. Pointer translation via the FOT. The pointer and the FOT are both contained in the same object (not shown). An FOT entry of 0 indicates an “internal” pointer.

Twizzler’s access control system prevents this from happening arbitrarily. The second system call is `invalidate_view`, which lets a thread inform the kernel of changed or deleted entries.

View objects not only reduce kernel boundary crossings, but they also improve the resumability of the system. After a power cycle, the OS now has information on which objects were mapped and where, improving the ability of threads to pick up where they left off. Additionally, view objects facilitate the sharing of address spaces between threads, since they can both synchronize on modifying a given view object and need not duplicate information. Note that the particular contents of a view object are system-specific. On virtual memory systems, one of their jobs is to manage ephemeral virtual mappings, while on other architectures their jobs may be to manage, e.g., segment tables. However, in all cases, views provide a mechanism for managing ephemeral state while providing enough context for threads to execute.

3.3 Persistent Pointers

Section 2 discussed the needs for references that outlive ephemeral actors. Twizzler provides *cross-object* persistent pointers so a pointer refers not to a virtual address but to an offset within an object by encoding an `object-id:offset` tuple. This enables a pointer to refer to persistent data, but it also allows objects to have *external* pointers that refer to data in any object in the global object space. We highlight cross-object pointers’ power and flexibility by demonstrating their ability to express inter-object relationships in Section 5.

To efficiently encode this tuple, we use indirection through a per-object **foreign object table (FOT)**, located at a known offset within each object. The FOT is an array of entries that each stores an object ID (or a name that resolves into an object ID, as we will see below) and flags. A cross-object pointer is stored as a 64 bit `FOT_idx:offset` value, where the `FOT_idx` is an index into the FOT. This provides us with both large offsets *and* large object IDs, since the IDs are not stored within the pointer itself. If an object wishes to point to data within itself (an *intra-object* pointer), then it stores 0 in `FOT_idx`. When dereferencing, Twizzler uses the `FOT_idx` part of the pointer as an index into the FOT, retrieving an object ID. The combination of a FOT and a cross-object pointer logically forms an `object-id:offset` pair, as shown in Figure 4.

Our design (discussed in prior work [10, 11]) differs from existing frameworks [6, 7, 14, 19, 20, 59] because of the indirection. Frameworks like PMDK store entire object IDs within pointers, increasing pointer size and reducing flexibility by removing the possibility of late-binding (discussed below). Additionally, Twizzler extends the namespace of data objects beyond one machine, as machine-independent data references are a natural consequence of cross-object pointers. Existing solutions are limited in this scalability. They either limit the ID space (necessary for storing IDs in pointers) and thus resort to complex coordination or serialization when sharing, or they require additional state (e.g., per-process or per-machine ID tables) that must be shared along with the data, forcing the receiving machine to “fix-up” references. Worse still, the fix-up is application-specific, since the object IDs are within any pointer, not in a generically known location. Our per-object FOT results in self-contained objects that are easier to share, thus interacting better with remote shared memory systems.

Part of our motivation for indirecting pointers through the FOT was to allow a large ID space without increasing pointer size. The density of NVM and the disaggregation of memory and applications means that we will be accessing data in a larger and larger address space, and it is vital that our abstractions allow for a large enough ID space to cover these needs. Since our IDs are 128 bits and our offsets need to support large objects, replacing pointers with a “fat pointer” style of just `object-id:offset` would mean more than doubling pointer size, which we found unacceptable. Other frameworks like PMDK, by contrast, increase pointer size to 128 bits for each pointer by encoding pointers as this tuple with 64 bit object IDs. The tradeoff is that our pointers take a little more work to translate (as they require an FOT lookup), but in return, we keep pointers 64 bits while supporting a truly global-scale address space. Thus, the overall space tradeoff is, for Twizzler, no additional space overhead per pointer, but an added 32-byte overhead per FOT entry. The number of FOT entries, however, is typically much smaller than the number of pointers, since pointers to the same external object can all use the same FOT entry. As we will see in Section 6, this has a dramatic benefit to performance.

FOT Entries and Late-binding. The FOT entry’s `flags` field has bits for read, write, and execute protections. The protections are *requests*; Twizzler implements separate access control on objects. This allows some pointers to refer to data with a read-only reference, while others can be used for writing, reducing stray writes (a single ID can repeat in the FOT with different protections). The FOT entries also enable atomic updates that apply to all pointers using that FOT entry.

Instead of *requiring* programmers to refer to objects via IDs only, we allow names in FOT entries. These entries may contain a pointer to an in-object string table that contains a name. Names enable late-binding [20], a vital aspect of systems, allowing references to objects that change over time, e.g., shared library versions. Names are passed to a *resolving* function (specified in the FOT entry). Allowing a program to specify how its names are resolved increases the flexibility of the system beyond supporting UNIX paths. Twizzler provides a default name resolver that uses UNIX-like paths.

The implementation of naming is orthogonal to Twizzler’s design. We allow a range of name resolution methods within the system stack and allow objects to specify their own name resolution functions for flexibility. For example, objects could be organized by both a relational database and a hierarchical namer similar to conventional file systems. Non-hierarchical file systems are well studied [4, 32, 33, 56, 57], but these systems do not easily cooperate atop a single data space. Since Twizzler uses a flat namespace as its “native” object naming scheme, it enables the required cooperation.

Pointer Translation. Current processors provide only a virtual memory abstraction, so applications must do some extra work to dereference a pointer, *translating* a pointer from its persistent form into a virtual address. This does not affect the *stored* pointer, which is still persistent and independent of any translation or address space. Thus, multiple applications, possibly with different address space layouts, can translate the same pointer at the same time without coordination.

Pointer translation occurs with the help of two `libtwz` functions: `ptr_lea` (load effective address) and `ptr_store`. When a program dereferences a pointer, it first calls `ptr_lea`. The pointer is resolved into an object-ID and offset pair through a lookup in the FOT, after which `libtwz` determines if the referenced object is already mapped (by maintaining per-view metadata). If not, then it picks an empty slot in the view and maps the object there (a cheap operation that does not invoke the kernel). Once mapped, `libtwz` combines the object’s temporary virtual base address with the offset and returns the new pointer. The `ptr_store` function does the opposite of `ptr_lea`—it turns a virtual pointer into a persistent one. While these are done manually in our implementation, we plan to implement compiler support to emit these calls automatically.

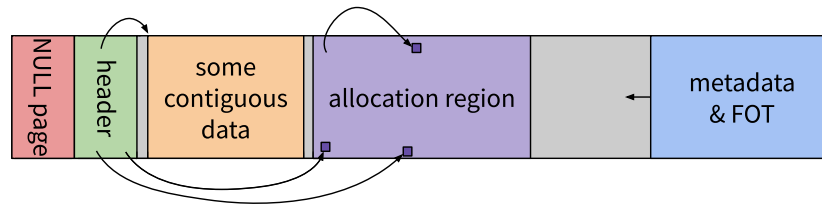


Fig. 5. A typical object layout. The header, contiguous region, and allocation region are all optional; however, most objects will have a header. This object contains a number of internal pointers between regions. The metadata region (which includes the FOT) grows downward.

FOT management is handled by `libtwz`. While a lookup in the FOT is a simple array-indexing operation, a store may require adding to the FOT. To avoid duplicate entries, `libtwz` walks the FOT looking for a compatible entry. If one is not found, then it atomically reserves a new entry and fills it (flushing cache-lines to persist it) before storing the pointer. The `ptr_store` operation is less common than `ptr_load`, and in the future, we may include additional caching metadata that would speed up the FOT walk (such as storing recent IDs).

Translating pointers has a small overhead (Section 6), and the result can be cached. Twizzler improves performance via a per-object cache of prior translations. The common case, intra-object pointers, does not require an external lookup and is implemented as a simple bitwise-or operation.

Object Allocation and Base Structures. Objects in Twizzler often have a header at the object’s base, the contents of which depend on what the object contains. Often these headers have pointers to other data in the object and describe the type of the object. For example, in our evaluation, we implement a red-black tree in an object. The header contains some basic information about the tree as well as a pointer to the root node. Placing headers at the object’s base gives applications a “starting point” that they can use to start accessing object data. Twizzler provides a dedicated function to get a pointer to an object’s header, called `obj_base`.

Note that the base address of an object is *not* at offset 0, but instead one page up, so we can still trap NULL pointers. If this were not the case, then a pointer value of 0 would still be a valid pointer, and we want to remain backwards-compatible with the assumption that a NULL pointer has integer value 0. The bottom page of an object is unmapped by Twizzler, allowing NULL pointer dereferences to be trapped by the kernel.

While objects are flat, contiguous regions of memory, different applications may want to organize that memory in different ways. Some objects, such as views, are largely interpreted as an array, but sometimes applications need to explicitly allocate and deallocate memory within an object. Twizzler provides an API to allocate and free units of memory from application-specified regions within objects. We make use of this in our red-black tree code, where new nodes are allocated out of the object using this API.

Figure 5 shows a typical object in Twizzler. The NULL page is always present to trap NULL pointers and is followed by a header. The application setting up this object may have a region of some contiguous data (such as some strings or an array) and may point to it from the header. The object may have a region setup for allocation so a future application using this object can easily allocate and free memory when manipulating the object. Finally, the FOT and metadata regions start at the top of the object and grow downwards.

3.4 Security and Access Control

Twizzler’s focus on memory-based objects requires that we design the security model around hardware-based enforcement, where the MMU checks each access. This design is *inevitable* in a data-centric OS, since the kernel is not involved in every memory access. The kernel merely

specifies the access rights when mapping an object and then relies on the hardware to enforce those rights with a low overhead.

A key design choice we make is *late-binding on security*. Applications request access to an object with permissions that they desire; if they access the object in only allowed ways (e.g., only reading a read-only object), then no fault occurs. This is because when we map an object (via a view), the kernel is not immediately involved, and so cannot check access rights for a particular access at the time the mapping is setup. Performing an access rights check on time of first access does not make sense either, as it associates a specific access (that might be allowed) with a permissions error. For example, if a program reads object *A*, and that program is allowed to read *A*, then it should be allowed to perform the read even if it requested read-write access to the object. This late-binding enables simpler programs that need not worry about elevating access rights through remapping data objects. Programs can make progress without knowing in advance the permissions of the objects they might access, thus enabling the reuse of the OS's access control mechanism in applications. We will show the flexibility of this in Section 5, wherein we add access control to a program by changing only a few lines of code.

Threads run in a security context [9, 26, 46], which contains a list of access rights for objects and allows the kernel to determine the access rights of programs. Using these contexts, Twizzler is able to provide analogues to groups and owners in UNIX while providing more fine-grained access control if necessary. Unlike past exploration into security contexts, data-centric OSes offer an advantage in simplicity. A security context abstraction in a UNIX-like OS needs to maintain access rights to a set of fundamentally different things (such as paths, virtual memory locations, and system calls). Instead, Twizzler's security contexts specify access rights to an object via IDs instead of virtual addresses. This also makes security contexts persistent, allowing us to use them as the primary way we assign security roles to threads. Security contexts are implemented via virtualization hardware (discussed in Section 4.1).

4 IMPLEMENTATION

Twizzler's kernel is similar to many microkernels, providing a small set of key primitives. It is 5,500 lines of architecture-independent code and 5,700 lines of architecture-dependent CPU driver code. The primary complexity in the system is implemented in userspace, as the design of the programming model greatly simplifies the kernel. Twizzler is open-source; more information can be found at <https://twizzler.io>.

4.1 Security Contexts and Page Tables

In Section 3.2, we discussed how view objects allow applications to specify what objects they want mapped in, and with what protections. These are merely *requests*, however; of course, if the thread does not have the appropriate *permissions*, then the kernel will program the address translation hardware appropriately. This presents a problem: Since threads can attach to a number of different security contexts, the number of different page-table structures that the kernel needs to manage grows quickly.

Twizzler uses Intel's **Extended Page Table (EPT)** technology,² which is part of the virtualization extensions. The EPT allows a virtual address to be translated by two separate page tables and is commonly used to virtualize the MMU in virtual machines. The first level, using normal MMU page tables, translates a virtual address to an object-logical address (typically with second-level address translation, this is referred to as the "guest-physical" address), and the second-level translates this to a physical address. This two-level translation scheme is shown in Figure 6.

²AMD has a similar system, but Twizzler does not support it yet.

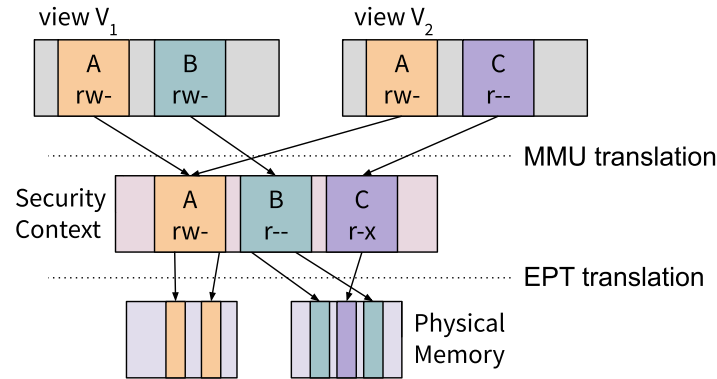


Fig. 6. Two-level translation scheme. The top level maps objects in virtual memory (defined by view objects) to those in object-logical space. The layout of the second level address space is managed by the kernel, since it is not user-facing, and the permissions are derived from a security context object. The second level maps to physical memory.

Two-level address translation via the EPT enables Twizzler to split protection *requests* and access control permissions. At the top level, Twizzler applies the requested protections from the view maps without restriction and programs the EPT to enforce access control derived from security context. Splitting protection and access control is what allows applications to simply map objects in for whatever access mode they would like without having to worry about first checking permissions. Furthermore, by separating out the permissions enforcement from ephemeral state and location mapping, we reduce the number of page table structures the kernel needs to manage from $O(nm)$ to $O(n + m)$, where n is the number of views and m is the number of security contexts. Views and security contexts can also be switched out independently from each other, which more closely fits the semantics of Twizzler.

Two-level mapping also greatly simplifies the design of the kernel. Since mapping objects to physical memory is done at the second level, page eviction is easy—the kernel can simply modify the shared page tables stored per-object, which updates the translation for all views and contexts on the system (after appropriate coherence, of course). Moving objects between DRAM and NVM is made easier, because objects reside in a given location within object-logical space regardless of where they are mapped in virtual memory, so the kernel does not need to maintain back pointers to update page table structures.

Virtualization Hardware. Twizzler’s use of virtualization hardware for normal operation is a limitation of existing processors. Intel does not have a mechanism for using the EPT without switching on the entire virtualization system and running in VMX-non-root mode. However, in practice, the additional overhead from running with virtualization is negligible, because we do not need all the protection of a traditional virtual machine and so we can switch much of it off. Because Twizzler’s kernel is its own guest, we can avoid much of the overhead introduced by VM exits necessary in lower-trust VM models. For example, Twizzler’s kernel is allowed to modify the EPT structures itself, despite being virtualized, and modern processors contain extensions that allow the guest to switch out EPTs itself and handle EPT faults without triggering a VM exit.

This pairs nicely with using the IOMMU as well—since EPT structures on Intel can be reused in the IOMMU, we can apply security contexts to drivers as well, making driver code less of a special case. For example, Twizzler provides a driver model for userspace drivers that allows driver code to construct security contexts that explicitly map in only the necessary objects that a device might need to access (e.g., command queues, data objects). As hardware devices grow in complexity and

increase their autonomy, treating them as additional computation resources and limiting access to objects through mechanisms already in place for normal applications allows simpler programming of advanced hardware devices.

4.2 Crash Consistency

Twizzler provides primitives for building crash-consistent data structures. At a low level, it provides mechanisms for writing back cache-lines, appropriate fences, and basic transactions. Applications use these primitives today outside of Twizzler to build up larger, more complex support for crash-consistent data structures.

Our goal is to provide low-level primitives without restricting programs or prematurely prescribing particular solutions. There is a wealth of research on crash-consistent data structures for NVM [16, 17, 25, 48, 52–55, 66], but it is still in flux. Of course, Twizzler manages *system* data structures, such as FOT entries, views, and so on, in a crash-consistent manner using the aforementioned primitives, locking, and fencing.

Twizzler also provides a transactional-persistent logging mechanism. Programmers can write TXSTART–TXEND blocks to denote transactions and TXRECORD statements to record pre-changed values. This is similar to the mechanism provided by PMDK [59]. If applications need more complex transactions using different logging mechanisms, then they can use libraries. Twizzler’s internal data structures and libtwz’s manipulation of object metadata is handled via a combination of these transactions and cache-line writebacks.

Twizzler provides a mechanism for restarting threads when power is restored following a crash. Since views are persistent objects, all objects mapped during a thread’s execution are known across power cycles and are mapped back in. The thread is then started at a special `_resume` entry point, allowing the program to handle the power failure in an application-specific manner. Of course, *volatile* objects will be lost when power resumes, and thus any attempted access to these objects will result in an exception. Thus, applications that wish to resume after power failure will need to be aware of and handle this. We do not wish to prescribe any restrictions here—applications that want to place their heap in volatile memory for performance or security reasons should be allowed to. We expect higher-level support for applications to manage persistent data, such as language support for persistent heaps, to make use of the features we provide, so applications that want to resume can put resuming information in persistent objects.

The reason we choose to restart threads at a known, different entry point from normal application startup is that in current systems, there is always volatile computation state (e.g., registers, the cache) that is lost when power is lost. Of course, in the future, systems may be able to prevent the loss of more and more ephemeral computation state (with the logical extreme being perfect resumability). In this case, the `_resume` handler can be a simple stub that resumes the execution exactly as left off. The more likely case, periodic checkpointing, can be similarly handled, with the `_resume` handler selecting the most recent valid checkpoint to resume from. The `_resume` handler enables all of these solutions, thus remaining applicable across hardware evolution.

4.3 Threading

Twizzler provides a set of threading primitives for applications. Threads in Twizzler are always attached to a view and one or more security contexts. Threads may communicate with each other using shared memory and can signal each other with a system call. Since everything in Twizzler is an object, each thread has a state object associated with it. Signals can be raised assuming the raiser has appropriate permissions on the state object, and the state object contains information about the thread.

A key primitive in Twizzler is the thread-sync system call. This call operates similar to `futex(2)` on Linux, except that it supports waiting on and waking up a number of different words of memory simultaneously. Multi-word thread-sync is necessary to support `select(2)`-like or `poll(2)`-like operations in a system where all data access is done with memory semantics. Twizzler’s standard library exposes an API for event handling that uses multi-world thread-sync, where objects may expose a set of “events” that can be triggered and waited for. This is used in numerous places to implement event handling for multiple communications streams implemented in objects.

4.4 FreeBSD Prototype

We also built a prototype of Twizzler by modifying the FreeBSD 11.0 kernel before implementing our standalone kernel. This was done both to more rapidly verify our design and to provide a prototyping environment for developers to write code for Twizzler in a familiar environment. We added Twizzler services to FreeBSD by adding system calls, modifying the fault-handling logic, and distinguishing Twizzler threads from FreeBSD threads. This is also a testament to the simplicity of the kernel in our model, since FreeBSD was relatively easy to modify to support the Twizzler userspace. However, the FreeBSD prototype is limited by its need to coordinate with FreeBSD’s UNIX services, thus the standalone kernel is more efficient and simpler and provides a better environment for researching kernel design changes in the face of NVM.

5 EVALUATION

Our primary goals for evaluating Twizzler were:

- (1) Show that Twizzler meets the needs of a data-centric OS in enabling programs to directly access persistent data.
- (2) Demonstrate that the programming model we defined provides sufficient power to easily and effectively build real applications with NVM in mind.
- (3) Measure the performance of our system to understand where we gain and lose performance.

We approached these goals two ways: porting existing software (SQLite) and writing new software for Twizzler. The first demonstrates both the generality of the programming environment (legacy software can be easily ported) and the potential performance gains to be had even for legacy software. The second demonstrates the true power of Twizzler’s programming model and allows us to explore the consequences of our design choices fully without being constrained by legacy designs.

We built three pieces of new software: a hash-table based **key-value store (KVS)**, a red-black tree data structure, and a logging daemon. Each had different characteristics and goals, and together they demonstrate the flexibility that Twizzler offers in allowing simple implementation, nearly free access control, and the ability to directly express complex relationships between objects. Using our KVS and red-black tree code, we ported SQLite (a widely used SQL implementation) to Twizzler along with a YCSB [18, 30] driver (a common benchmark), allowing us to explore Twizzler’s model in a larger, existing program that would let us study the performance of Twizzler in a complex system that stores *and processes* data. We present the performance of SQLite and our new software, along with microbenchmarks, in Section 6.

5.1 Case Study: Key-value Store

We implemented a multi-threaded hash-table based **key-value store (KVS)**, called `twzkv`, to study cross-object pointers and our late-binding of access control. Our KVS supports insert, lookup, and delete of values by key (both of arbitrary size) and hands out direct pointers to persistent data

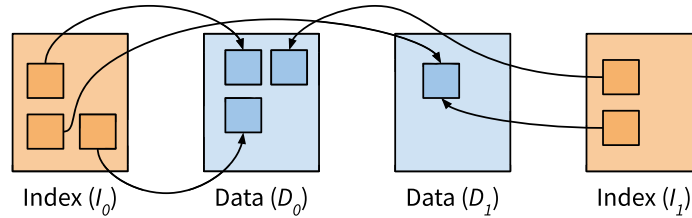


Fig. 7. Cross-object pointers in twzkv. The index object contains pointers to keys and values, which can reside in any data object. We also support multiple indexes to enable additional indexing strategies and access control on discovery. Because Twizzler provides native support for cross-object pointers, this kind of construction is no more complex than keeping all keys and data within a single object.

during lookup. During insert, it copies data into a data region before indexing the inserted key and value. We built twzkv in multiple phases to study how our system handles changing requirements.

We built twzkv in roughly 250 lines of C. Handing out direct pointers into data was trivial to implement with cross-object pointers, requiring only a call to `ptr_lea` during lookup. The initial implementation maintains two objects, one for data and one for the index. The complexity typically involved when storing both index and data in a single, flat file is not justified in a programming model where we can express inter-object relationships directly at near-zero cost in complexity or performance. In our case, a pointer from the index object to the data object (such as an entry in the hash table) can be written with a single call to `ptr_store`. This, combined with the simple requirements for an in-memory NVM KVS, resulted in a small implementation that was nonetheless a usable KVS.

Extending Requirements. Next, we added functionality to protect values with access control. We wanted to keep handing out direct pointers to data during lookup and to keep twzkv a library (as opposed to a service). Meeting these goals on an existing system would be difficult without adding significant complexity, such as reimplementing a lot of Twizzler’s pointer framework or implementing manual, redundant access control.

In Twizzler, implementing access control in twzkv involved having the index refer to data in multiple data objects, assigning those objects different access rights, and allocating from those objects depending on desired access rights. We were able to implement this while preserving the original code due to the transparent nature of Twizzler’s cross-object pointers. Now, when inserting, the application indicates the data object into which to copy the data, as shown in Figure 7.

By supporting multiple data objects, twzkv can leverage the OS’s access control, minimizing complexity. Unrestricted data can go in D_0 (Figure 7), whereas restricted data can go in D_1 . Since each object has distinct access control, a user can set the objects’ access rights, then decide where to insert data according to policy. The indexes point to the correct locations regardless of the access restrictions of the data objects, and twzkv still hands out direct pointers, but a user that is restricted from accessing data in D_1 will not be able to dereference the pointer. A further extension is to support secondary indices, as shown in Figure 7, enabling alternative lookup methods and limiting data discovery with index object access control. This extension is easy to implement on Twizzler.

Comparison to UNIX Implementation. To compare with existing techniques, we built a similar KVS using only UNIX features (called `unixkv`). It also separates index and data, but it must manually compute and construct pointers, requiring a significant amount of programmer time to get right. Supporting multiple data objects was complex in `unixkv`, because we had to store and process file paths in the index and store references to paths for pointers, increasing overhead and code complexity by 36%—a lot for an implementation with relatively few pointers—just to reimplement

Twizzler’s support. The extra complexity also included code to manually open, map, and grow files, much of which Twizzler handles internally. Development time was extended by bugs that were not present when developing twzkv, due to the manual pointer processing. While twzkv gains transparent access control, unixkv does not due to the lack of on-demand object mapping and late-binding of security. Instead, unixkv needs to know object permissions before mapping, a restriction that limits the ability to reuse OS access control, something that twzkv could leverage through late-binding on security (Section 3.4).³ Other frameworks like PMDK that do not integrate access control and late-binding into their models have similar limitations.

5.2 Case Study: Red-black Tree

To evaluate the process of writing persistent, “pointer-heavy” data structures, we implemented a red-black tree in C using normal pointers (ramrbt) in 100 lines of code and evolved it for persistent memory in two ways: manually writing base+offset style pointers, as current systems require (unixrbt), and using Twizzler (twzrbt). Porting existing data structure code to persistent memory will be common during the adoption of NVM, and much of the complexity therein comes from dealing with persisting virtual addresses [49].

In developing unixrbt, we found 83 locations where we had to perform pointer arithmetic for converting between persistent and virtual addresses. Consider an expression such as `root->left->right = foo`. Inserting calls to translate this directly results in `L(L(root)->left)->right = C(foo)`, where `L` converts to a virtual address and `C` converts back, which is heavily obfuscated and took more development time than writing ramrbt in the first place due to debugging.

We built twzrbt like unixrbt, annotating pointer stores and dereferences. However, unixrbt used an application-specific solution for pointer management; if other applications wanted to use the data structures created by unixrbt, they would have to know the implementation details of the pointer system (or share the implementation, thus reimplementing much of Twizzler’s library). Additionally, due to Twizzler enabling improved system-wide support for cross-object pointers, these transformations can be made automatic through compiler and linker support.

Unlike twzrbt, unixrbt’s tree is limited to a single persistent object; a limitation that prevents the tree from growing arbitrarily, does not allow it to directly encode references to data outside the tree object, and does not gain it the benefits of cross-object data references that were discussed above for twzkv. Adding support for this to unixrbt would require modifying the core data structures to include paths and significantly altering the code, increasing its length by at least a factor of 2, whereas twzrbt gets this functionality for free.

Another advantage of twzrbt is reduced support code compared to unixrbt; unixrbt needed code to manage and grow files and mappings, while we implemented twzrbt as simple data structure code with Twizzler managing that complexity. The additional error handling code and pointer validity checks in unixrbt (handled automatically in Twizzler) increased development time and implementation complexity.

5.3 Porting SQLite

We ported SQLite to Twizzler to demonstrate our support for existing software and to evaluate the performance of a SQLite backend designed for Twizzler. We used our POSIX support framework, a combination of musl and our library twix, to support much of SQLite’s POSIX use. We took a modified version of SQLite called SQLLightning that replaced SQLite’s storage backend with a memory-mapped KVS called LMDB [15]. We chose this port because LMDB is implemented with

³unixkv could trap segmentation faults to do this, but that would be application-specific, difficult, and would reimplement Twizzler functionality.

mmap'd files as the primary access method and hands out direct pointers to data as one would expect from an effectively designed NVM KVS.⁴ Since LMDB's SQLightning port already replaces the storage backend with calls to LMDB, we ported SQLite to Twizzler by taking our KVS and red-black tree code and implementing enough of the LMDB interface for SQLite to run using Twizzler as a backend. Outside of the B-tree source file few changes were needed for SQLite to run on Twizzler. We further ported our modified SQLite backend to PMDK to compare directly with a commonly used NVM programming library that supports persistent pointers.

We also ported a C++ YCSB driver [30], which required porting the C++ **standard template library (STL)**. Since we had already ported a standard C library, the C++ STL was easily ported, demonstrating the ease of porting software to Twizzler. We have also ported some existing UNIX utilities (such as bash and busybox), which largely require only recompiling to run on Twizzler. Of course, to gain *all* of the benefits of Twizzler, programs will need to be written with NVM in mind (but this is true regardless of the target OS).

Our implementation of the LMDB interface corroborated our experience from the KVS case study: Much of the complexity in storage interfaces and implementations comes from the separation between storage and memory. This has been studied before (as we will elucidate in Section 7), but the advent of NVM changes the game significantly by allowing programmers to think directly via in-memory data structures. The result is that interfaces like cursors in a KVS become redundant. We implemented this interface for LMDB, but the functions were largely wrappers around storing a pointer to a B-tree node and traversing the tree directly without separate loads and copies. The result was an extremely simple implementation (500 LoC) that still met the required interface. Future software for NVM can use Twizzler's programming model to more effectively write software that eschews the need for complexity forced by the two-tier storage hierarchy.

5.4 Porting Additional Applications

In general, porting in Twizzler is straightforward. We have a collection of tools that provide a framework for compiling software using the Twizzler toolchain against other ported software and libraries. Since we have chosen musl as our standard C library, many applications work already with minor changes. However, it is often the case that applications require some small tweaks to get running—for example, configuration paths—an experience common for anyone who has ported software to a new operating system.

To date, we have ported a number of tools one would expect to find on a UNIX system, such as busybox (providing numerous command-line utilities), bash, vim, gcc, binutils, and others. Many of these programs required little or no modification. Of course, this means that they do not gain some of the benefits Twizzler's model provides, since they still operate on persistent data with a POSIX model, however our goal in porting these tools was not to improve their performance, it was to provide a somewhat familiar environment for users.

Of course, perfect emulation of a Linux kernel is a huge effort, and it is not the primary goal of our research. As a result, not all system calls are implemented and Linux features like `procfs` are lacking. This means that some programs may require features that are not yet implemented, and therefore require modifications to `twix` to run. However, as we continue to port software, `twix`'s coverage of Linux features grows, making future porting easier. We will continue to implement more support in `twix` for applications as needs arise. Note that many applications (even complex applications like `gcc`) often boil down to reading and writing files and managing processes, all of which is implemented.

⁴These are not persistent pointers, however, unlike Twizzler's.

5.5 Discussion

Although these implementations were simple, they represent the applications and data structures we expect in a data-centric system. Persistent pointers we can directly use in our programming languages make computing over persistent data almost transparent, allowing simple implementations that are nevertheless easy to evolve as requirements change.

Not only does twzkv have strong support for access control, it enables concurrent use of databases via cross-object pointers. Applications can load indexes for multiple databases without needing to worry about address space layout and without writing complex pointer management code that would be required by an implementation using `mmap`. We were able to provide access control without a single line of code in twzkv dedicated to checking or enforcing access rights. Instead, we relied on Twizzler’s built-in access control, something not possible with other frameworks that do not support late-binding of access rights and do not consider security as part of their programming model. Twizzler thus removes the need for applications to enforce and implement their own access control, which increases the security of the system by divesting programmers from the responsibility of getting the enforcement right. Similar functionality for current systems would traditionally require separation of the library and application into a client-server model, but that additional overhead is unneeded here and inappropriate on a persistent memory system.

Although twzrbt and twzkv had different densities of pointer operations, twzrbt being “pointer-heavy” and twzkv being “pointer-light,” Twizzler improved the complexity of both over manual implementation and improved flexibility over existing persistent pointer methods. Using a system-wide standardized approach to pointer translations not only enables better compiler and hardware support, but it also improves interoperability; because they share a common framework, twzkv could use the red-black tree code and data with ease, and even interact with the SQLite database even though they were written separately without that goal in mind. The position-independence afforded by this model enables both composability and concurrency, while also simplifying programming on persistent data to a natural expression of data structures.

Non-shared-memory Programs. To push the limits of our model and show that Twizzler does not constrain programmers into a shared-memory model, we implemented a logging framework (similar to `syslogd`). The logging daemon, `logboi`, receives logging requests through a shared stream object (an API provided by Twizzler). This connection is setup via our security mechanism that allows for threads to switch security contexts, enabling the construction of secure calls to set up private communication channels.

Once a private channel is set up, the application can stream logging events to `logboi`, which collects them as any logging daemon would. The implementation of `logboi` must be able to handle multiple clients and clients crashing unexpectedly. This demonstrates the flexibility of our model: Despite being completely in userspace, this communication stream implements a message-passing style model. When handling multiple clients, `logboi` uses the multi-word thread-sync system call that we discussed earlier to wait for one of a number of possible events. When a thread exits, the kernel updates the information in the thread’s state object and performs a thread-sync wakeup on that word, allowing `logboi` to detect when a thread has exited even if it does so unexpectedly.

6 PERFORMANCE

Our evaluation’s primary focus is on the benefits of the programming model, showing new functionality with reduced complexity at an acceptable overhead. Nevertheless, there are many cases where we see significant improvement (such as SQLite), because the programming model has less overhead, and our pointer design is spac-efficient and fast to translate.

Table 1. Latency of Common Twizzler Operations, Including Pointer Loading and Storing, and Object Mapping

Pointer Resolution Action	Average Latency (ns)
Uncached FOT translation	27.9 ± 0.1
Cached FOT translation	3.2 ± 0.1
Intra-object translation	0.4 ± 0.1
Inter-object pointer store	17.2 ± 0.6
Intra-object pointer store	2.3 ± 0.1
Mapping object overhead	49.4 ± 0.2

We measured the performance of our KVS and red-black tree, performed microbenchmarks, and evaluated the Twizzler port of SQLite against Linux (Ubuntu 19.10) instances of SQLite, SQLite-Lightning, and our port of SQLite to PMDK. Tests ran on an Intel Xeon Gold 5218 CPU running at 2.30 GHz with 192 GB of DRAM and 128 GB of Intel Persistent DIMMs. We compiled all tests against the `musl` C library instead of `glibc`, because Twizzler uses `musl` to support UNIX programs.

All Linux tests used the NOVA filesystem [69] (a filesystem optimized for NVM) on the NVDIMMs, mounted in DAX mode. This enabled direct access to the persistent memory without a page-cache interposing on accesses.

6.1 Microbenchmarks

Table 1 shows common Twizzler functions’ latencies, including pointer translation (loading and storing) and mapping overhead. The overhead shown for resolving pointers does not include dereferencing the final result, since that is required regardless of how a pointer is resolved. The first row shows the latency for resolving pointers to objects the first time. Twizzler makes a further optimization by caching the results of translations for a given FOT entry. Each successive time that FOT entry is used to resolve a pointer, the result of the original translation is returned immediately, improving the latency as shown on the “cached” row of Table 1. Note that the low latency of these results is expected; the performance critical case of these functions’ use is repeated calls, and, since these operations are simple, they fit within the processor cache.

Twizzler translates intra-object pointers by first checking if the pointer is internal and, if so, adding the object’s base address to it—the same operation required for application-specific persistent pointers. The expanded programming model offered by Twizzler makes this overhead minor relative to the high costs for persistent data access on current systems, which have high-latency for equivalent operations.

The pointer store operations shown in Table 1 measure the latency of the `ptr_store` operation that is used to construct persistent data references in Twizzler. While these operations are less common than pointer loads, their overhead directly affects applications that perform many updates to data structures. The most common pointer store operation applications perform is internal (intra-object) pointer stores, in which the overhead is minimal. Pointer store operations for external (inter-object) references have slightly more overhead, since they need to perform an FOT scan to allow FOT entry reuse.

We compared our pointer translation to UNIX functions. Resolving an external pointer with an ID corresponds roughly to a call to `open(‘‘id’’)`, which has a latency of $1,036 \pm 15$ ns. The comparison is not exact, of course; the pointer resolution also maps objects, and the call to `open` must handle file system semantics. However, the direct-access nature of NVM results in pointer translation achieving the same goal as opening a file does today. The pointer operations in Twizzler

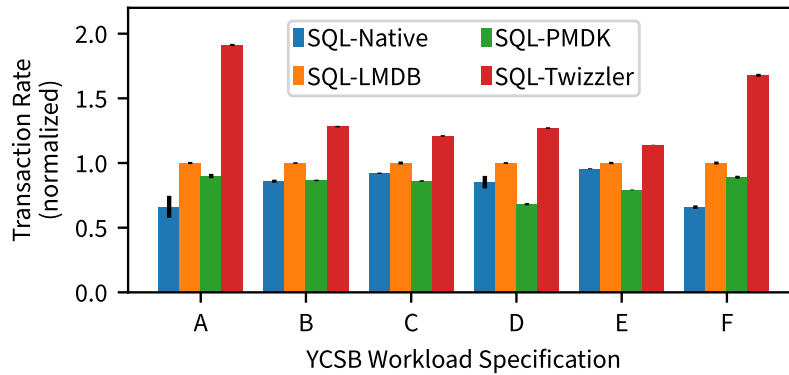


Fig. 8. YCSB throughput, normalized (higher is better). Twizzler outperforms all other variants in all tests.

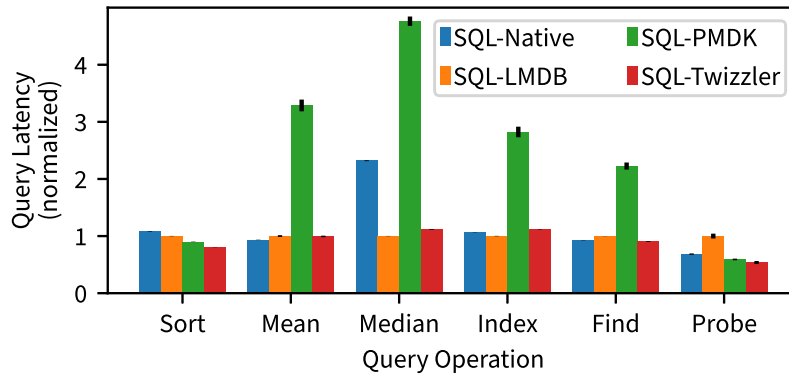


Fig. 9. Query latency, normalized (lower is better). Twizzler maintains a similar level of performance with the native and LMDB variants, despite comparing Twizzler’s simplistic red-black tree index implementation with highly optimized B-trees. Twizzler also significantly outperforms PMDK despite sharing a similar implementation for the index structures.

accomplish much of the same functionality as the heavier-weight I/O system calls on UNIX with more utility and less overhead.

A more direct comparison is object mapping, which has low latency compared to `mmap` (658.7 ± 12.7 ns—a $13.3\times$ speedup) though the two have similar functionality. Since mapping occurs entirely in userspace, cache pollution is reduced. While both `mmap` and Twizzler’s mapping require page-faults to occur before the data is actually mapped, this overhead is similar in Twizzler and UNIX, and so is not shown.

6.2 SQLite

We ran four variants of SQLite—three on Linux and one on Twizzler—and compared their performance: “SQL-Native” (unmodified SQLite), “SQL-LMDB” (SQLite using LMDB as the storage backend), “SQL-PMDK” (SQLite using our red-black tree on PMDK), and “SQL-Twizzler” (our port of SQLite running on Twizzler). SQL-Native was run in `mmap` mode so both it and SQL-LMDB used `mmap` to access data. We ran each on the same hardware and normalized the results.

Figure 8 shows the three variants’ throughput under standard YCSB workloads. The performance improvement of the LMDB and Twizzler variants over SQL-Native is likely due to handing SQLite direct pointers to data. However, in the Twizzler case, we get an additional benefit of operating on data structures directly while LMDB has an abstraction cost.

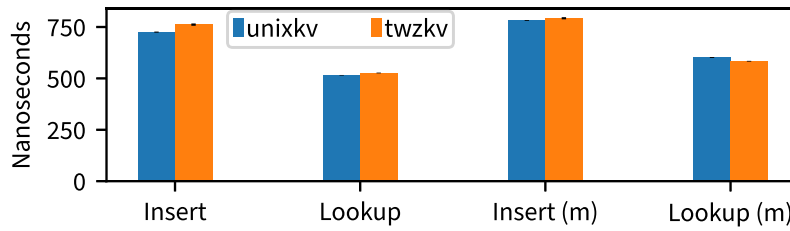


Fig. 10. Latency of insert and lookup in twzkv and unixkv. An “(m)” indicates support for multiple data objects. Both unixkv and twzkv have similar latencies.

Figure 9 shows the latency of queries on a 1M row table. This is common data processing—loading and then examining data in a variety of ways. We measured the performance of calculating the mean and median, sorting rows, finding a specific row, building an index, and probing the index. SQL-Twizzler had similar performance to SQL-LMDB and SQL-Native despite comparing its extremely simple storage backend to optimized B-tree backends (that benefit from scan operations). As a more direct comparison, SQL-Twizzler significantly out-performed SQL-PMDK in most tests. PMDK’s pointer operations are more expensive than Twizzler’s, requiring up to two hash table lookups per translation [6]. Additionally, PMDK’s pointers are 128 bits, while Twizzler does not increase pointer size. Increased pointer size results in significantly worse cache performance, especially in a pointer-heavy data structure like a persistent red-black tree.

6.3 Key Value Store

We compared twzkv to unixkv by inserting 1M distinct key-value pairs, followed by looking up each in order. The inserted items were 32-bit keys and 32-bit values, chosen to reduce the overhead of data copying, since we were focusing on pointer translation overhead. Both were compared under two modes: single-data-object and multiple-data-objects. Both KVSes translated between virtual and persistent addresses when storing and retrieving data, but for multiple-data-objects, we allow for storing the data in an arbitrary object.

Figure 10 shows the latency of lookup and insert, demonstrating that not only is the memory-based index and data object structure that can hand out direct data pointers sufficiently low latency to take advantage of NVM, but the additional overhead of cross-object pointers is minimal. Compared to unixkv, twzkv has minimal overhead in the single-object case and improves lookup performance in the multiple-object case. The minor overhead in other cases comes with improved flexibility, simplicity, and access control support (unixkv does not support access control). Finally, multithreaded access on twzkv and unixkv did not improve performance; despite the pointer translations, they ran at memory bandwidth (for NVM).

6.4 Red-black Tree

We measured the latency of insert and lookup of 1M 32-bit integers on both unixrbt and twzrbt. The insert and lookup latency of twzrbt was 528 ± 3 ns and 251.8 ± 0.5 ns, while insert and lookup latency of unixrbt was 515 ± 2 ns and 213 ± 1 ns. The modest overhead comes with significantly improved flexibility, as unixrbt does not support cross-object trees, and less support code (unixrbt manually implements mapping and pointer translations). Note that even though there is lookup overhead in twzrbt, this overhead did not predict the results of a larger program—the SQL-Twizzler port used this red-black tree and saw performance benefits over block-based implementations.

Table 2. Latency of Selected twix System Calls Compared to Linux System Calls

System Call	OS	Average Latency (ns)
getpid	Linux	98.7 ± 2.3
	Twizzler	10.2 ± 0.2
read	Linux	321.4 ± 0.2
	Twizzler	55.4 ± 0.2

6.5 Twix System Call Overhead

Our UNIX emulation layer, *twix*, is meant to provide compatibility for legacy applications. While we expect that applications will wish to take full advantage of NVM and Twizzler’s improvements in programmability and performance, we can still provide a small benefit for applications that rely on *twix* to provide POSIX-like I/O. Access to *twix* is done by *musl*, the C library we use, when it would normally perform a system call to a Linux kernel. We replaced all instances of the `syscall` instruction in C and assembly code in *musl* with a `call` instruction to an entry point in *twix*. This entry point, despite being a function call, obeys the Linux system call ABI (e.g., which registers hold parameters). Thus, while it has significantly less overhead than a full system call and context switch, it does still have higher overhead than a normal function call, since it must back up and restore all registers.

Table 2 shows the latency of some selected system calls on both Linux and Twizzler (implemented via *twix*). As expected, `getpid`’s overhead is small on both systems, but on Twizzler it is significantly lower. The difference, in this case, comes largely from the kernel entry overhead. A small amount of additional overhead comes from *twix* matching the Linux system call ABI and having to call its `getpid` implementation through a lookup table.

We also measured the latency of a call to `read` for a file. We chose to do reads on cached files for a small number of (already cached) bytes to avoid device transfer overhead. Performing a file read on Twizzler often amounts to a call to `memcpy`, so applications that perform large numbers of small reads could see some benefit. In contrast, on Linux, the kernel needs to traverse internal file structures, the page-cache, and possibly file system structures. However, as we said, *twix* is intended for legacy support, not performance improvement, despite the lower system call overhead.

7 RELATED WORK

Twizzler’s design is shaped by fundamental OS research [13, 19, 27–29, 43, 44], which, while approaching similar topics described in Section 2, often did not consider *both* design requirements simultaneously, resulting in an incomplete picture for NVM. Recent research on building NVM data structures [16, 17, 23, 39, 47, 66] often focuses on building data structures that provide failure atomicity and consistency. In contrast, we explore how NVM affects programming models while potentially improving performance. We draw from recent work on providing OS support for NVM systems [12] and work providing recommendations for NVM systems [50], integrating object-oriented techniques and simplified kernel design to provide high-performance OS support for applications running on a single-level store [5, 62].

Memory and Object Model. Multics was one of the first systems to use segments to partition memory and support relocation [7, 20]. It used segments to support location independence, but still stored them in a file system, requiring manual linkage rather than the automated linkage in Twizzler. Nonetheless, Multics demonstrated that the use of segmenting for memory management can be a viable approach, though its symbolic addresses were slow.

The core of Twizzler’s object space design uses concepts from Opal [13], which used a single virtual address space for all processes on a system, making it easier to share data between programs. However, Opal was a single-address space OS, which is insufficient for NVM [10, 11], and, while it resulted in a speedup of data transfer and sharing as well as interfacing with devices, it did not address issues of file storage and name resolution. It also still required a file system, since there was no way to have a pointer refer to an object with changing identity, whereas our approach of late-binding for pointers removes the need for an explicit file system. Other single-address space OSes, such as Mungi [36], Nemesis [58], and Sombrero [64], show that single address spaces have merit, but, like Opal, did not consider how the use of NVM would alter their design choices; in particular, how the use of fixed addresses results in a great deal of coordination that is unnecessary in our approach. OSes such as HYDRA [68] provide functionality similar to cross-object pointers; however, in Twizzler, we extend their use from procedures-referencing-data to a more general approach. Furthermore, they required heavy kernel involvement, an approach incompatible with our design goals.

Single-level stores [22, 61, 63] remove the memory versus persistent storage distinction, using a single model for data at all levels. While well-known, “little has appeared about them in the public literature” [61], even since the EROS paper. Our work is partially inspired by Grasshopper [22], AS/400, and orthogonal persistence systems, but while these are designed to provide an illusion of persistent memory, Twizzler is built for real NVM and focuses on providing a truly global object space with global references without cross-machine coordination. Clouds [21] implemented a distributed object store in which objects contained code, persistent data, and both volatile and persistent heaps. Our approach uses lighter-weight objects, allowing direct access to objects from outside, unlike Clouds. Software persistent memory [34], designed to operate within the constraints of existing systems, built a persistent pointer system using explicit serialization without cross-object references, in contrast to Twizzler. Meza [51] suggested hardware manage a hybrid persistent-volatile store with fine-grained movement to and from persistent storage. Since persistence in Twizzler is to NVM, we need not interpose on movement between storage and memory, instead simply managing memory mappings of persistent objects, reducing OS overhead.

Recently, several projects have considered the impact of non-volatile memories on OS structure. Bailey et al. [5] suggest a single-level store design. Faraboschi et al. [31] discuss challenges and inevitable system organization arising from large NVM, and we follow many of their recommendations. The Moneta project [12] noted that removing the heavyweight OS stack dramatically improved performance. While Moneta focused on I/O performance, not on rethinking the system stack, we leverage their approach to reduce OS overhead as much as possible, even when the OS must intervene. Lee and Won [45] considered the impact of NVM on system initialization by addressing the issue of system boot as a way to restore the system to a known state; we may need to include similar techniques to address the problem of system corruption. Our work evolved from some earlier work where we laid the groundwork for abstraction requirements for both hardware and software for NVM [10] and a discussion on the implications of system-wide persistent data references [11].

Object Model. IBM’s K42 [44] inspired the high-level design of Twizzler. The object-oriented approach to designing a micro or exokernel used in K42 is an efficient design for implementing modular OS components. Like K42, Twizzler lazily maps in only the resources that an application *needs* to execute. Similar techniques for faulting-in objects at runtime have been studied [38]. Communication between objects in Twizzler is, in part, implemented as protected calls, similar to K42.

Emerald [41, 42] and Mesos [37] implemented networked object mobility, which we can also support. Emerald implemented a kernel, language, and compiler to allow objects mobility using

wrapper data structures to track metadata and presenting objects in an object-oriented language, impacting performance via added indirection for even simple operations.

The Twizzler object model was shaped by NV-heaps [16], which provides memory-safe persistent objects suitable for NVM and describes safety pitfalls in providing direct access to NVM. While they have language primitives to enable persistent structures, Twizzler provides a lower-level and uninhibited view of objects like Mnemosyne [66], allowing more powerful programs to be built. Languages and libraries may impose further restrictions on NVM use, but Twizzler itself does not. Furthermore, Twizzler’s cross-object pointers allow external data references by code, whereas NV-heap’s and DSPM’s [60] pointers are only internal. Existing work beyond Multics on external references shows and recommends hardware support [59, 67], but provides a static or per-process view of objects, unlike Twizzler, limiting scalability and flexibility.

Projects such as PMFS [25] and NOVA [69] provide a file system for NVM. Twizzler, in contrast, provides direct NVM access atop of a key-value interface of objects. Although Twizzler does not supply a file system, one can be built atop it. While NOVA and PMFS provide direct access to NVM, NOVA adds indirection with copies. Both use `mmap` (which falls short as discussed above) and, unlike Twizzler, require significant kernel interaction when using persistent memory.

Our kernel that “gets out of the way” is influenced by systems such as Exokernel [29] and SPIN [8], both of which drew on Mach [3]. In Exokernel, much of the OS is implemented in userspace, with the kernel providing only resource protection. Our approach is similar in some respects, but goes further in providing a single unified namespace for all objects, making it simpler to develop programs that can leverage NVM to make their state persistent. In contrast, SPIN used type-safe languages to provide protection and extensibility; our approach cannot rely upon language-provided type safety, since we want to provide a general purpose platform.

8 FUTURE WORK

Compiler and Hardware Support. Twizzler’s clean-slate NVM abstraction reopens the possibility of coevolving OSes, compilers and languages, and hardware. Standardized OS support for cross-object pointers provides a stationary target for both compilers and hardware to design towards, whereas application-specific solutions do not. Twizzler’s pointer translation functions are simple enough to be emitted by a compiler. We plan to explore adding basic compiler support for C and C++ to automatically interoperate with Twizzler so persistent pointers are even more transparent to the compiler. Better still, we would like to study additional language-level support for persistent pointers, including type and lifetime annotations (such as the ones supported in Rust [2]) for additional semantics the compiler can make use of when emitting code that operates directly on persistent data structures.

Hardware support, too, can be helpful in improving the performance of our pointer translations. With Twizzler providing a common framework, we can clearly state our needs to hardware. For example, hardware-accelerated FOT access would improve the performance in pointer-heavy data structures. Segmentation support, allowing us to assign page-tables for each object and load them in as needed, would dramatically speed up memory mapping (and move memory management closer to the semantics of our programming model). Finally, first-class support for abstracting physical memory—a necessary feature for efficiently moving data around in a heterogeneous memory hierarchy in the face of numerous devices—would simplify the design of the kernel, because we would not need to invoke the entirety of the virtualization hardware. We are interested in exploring modifications to RISC-V to better support Twizzler.

Security. Although we discussed the Twizzler security model briefly, there is still much to do. The current model provides access control, a basic ability to define and assign roles based on security contexts, and simple sub-process fault isolation through the ability to switch security contexts.

We are exploring a *flexible* security model that allows programmers to easily trade off between security, transparency, and performance using capability-based verification. For example, we are implementing a call-gating mechanism that will allow us to restrict control-flow transfers between application components, improving the security against malicious components and reducing the possibility of memory-corrupting bugs.

Networking and Distributed Twizzler. One of the key principles of Twizzler is to focus the programming model on data and away from ephemeral actors such as processes and nodes. This is enabled by our identity-based pointers that decouple location from references, and by ensuring all the context necessary to understand these relationships is stored with the data. Because our data relationships are independent of the context of a particular machine, applications can more easily share data. This easy sharing, combined with a large ID address space, motivates a *truly* global object ID space.

We are building a networking stack and support for a distributed object space into Twizzler. Our networking stack is based around extensive use of hardware virtualization in modern NICs. This design, which is in use in existing kernel-bypass strategies, will work well with our core OS design of reducing kernel interposition. At a higher level, we are considering how distributed applications change in our model. For example, an increase in data mobility facilitated by our location-independent data references and identities means that we can manifest both data and code where they are needed without complex marshalling, turning distributed computation into a rendezvous problem. We plan to build distributed applications atop Twizzler to demonstrate this approach.

Of course, for compatibility, we will provide a traditional sockets-based networking stack. However, we can use existing userspace libraries that, e.g., implement TCP in userspace. Because we implemented our POSIX compatibility library in userspace, applications can gain many benefits afforded by kernel-bypass networking frameworks while still using traditional socket interfaces.

Alternative Block Storage Technologies. Our work meshes well with key-value SSDs, which extend the NVMe specification to include put and get operations. This would allow us to store and retrieve parts of objects based on their names rather than block addresses, thus greatly simplifying the storage system of the OS, because it removes the need for a filesystem. Twizzler uses a userspace pager design for moving data between memory and indirectly accessible storage; providing a more “native” interface for object-based storage will greatly improve the performance of this system. We have prototype KVSSDs, and are in the process of implementing support for them in Twizzler.

9 CONCLUSION

Operating systems must evolve to support future trends in memory hierarchy organization. Failing to evolve will relegate new technology to outdated access models, preventing it from reaching full potential and making it difficult for OSes to evolve in the future. Twizzler shows a way forward: an OS designed around NVM that provides new, efficient, and easy-to-use semantics for direct access to memory. Cross-object pointers in Twizzler allow programmers to easily build composable and extensible applications with low overhead by removing the kernel from persistent data access paths, thereby improving the flexibility and performance. Our simpler programming model improved performance despite the (small) pointer translation overhead. Even a memory hierarchy with large RAM but without persistent memory benefits from our design by enabling programs to operate on large, shared, in-memory data with ease. Our programming model is easy to work with compared to existing systems, and we were able to both quickly prototype real applications with advanced access control features and port existing software (SQLite). Twizzler

will give us a system from which we can build a full NVM-based OS around a data-centric design and explore the future of applications, OSES, and processor design on a new memory hierarchy.

Availability. Twizzler is available at twizzler.io.

ACKNOWLEDGMENTS

We would like to thank the members of the Storage Systems Research Center for their support and feedback.

REFERENCES

- [1] [n.d.]. The musl C Library. Retrieved from <https://musl.libc.org/>.
- [2] [n.d.]. The Rust Programming Language. Retrieved from <https://www.rust-lang.org/>.
- [3] Mike Accetta, Robert Baron, William Bolosky, David Golub, Richard Rashid, Avadis Tevanian, and Michael Young. 1986. Mach: A new kernel foundation for UNIX development. In *Proceedings of the USENIX Technical Conference*. USENIX, 93–112. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/accetta-usenix86s.pdf>.
- [4] Sasha Ames, Nikhil Bobb, Kevin M. Greenan, Owen S. Hofmann, Mark W. Storer, Carlos Maltzahn, Ethan L. Miller, and Scott A. Brandt. 2006. LiFS: An attribute-rich file system for storage class memories. In *Proceedings of the 23rd IEEE/14th NASA Goddard Conference on Mass Storage Systems and Technologies*. IEEE. Retrieved from <http://www.ssrc.ucsc.edu/Papers/ames-mss06.pdf>.
- [5] Katelin Bailey, Luis Ceze, Steven D. Gribble, and Henry M. Levy. 2011. Operating system implications of fast, cheap, non-volatile memory. In *Proceedings of the 13th Workshop on Hot Topics in Operating Systems (HotOS'11)*. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/bailey-hotos11.pdf>.
- [6] Piotr Balcer. 2015. An introduction to pmemobj (part 1) - accessing the persistent memory. Retrieved from <https://pmem.io/2015/06/13/accessing-pmem.html>.
- [7] A. Bensoussan, C. T. Clingen, and R. C. Daley. 1969. The Multics virtual memory: Concepts and design. In *Proceedings of the 2nd ACM Symposium on Operating Systems Principles (SOSP'69)*.
- [8] Brian N. Bershad, Stefan Savage, Przemyslaw Pardyak, Emin Gün Sirer, Marc E. Fiuczynski, David Becker, Craig Chambers, and Susan Eggers. 1995. Extensibility, safety, and performance in the SPIN operating system. In *Proceedings of the 15th ACM Symposium on Operating Systems Principles (SOSP'95)*. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/bershad-sosp95.pdf>.
- [9] Andrea Bittau, Petr Marchenko, Mark Handley, and Brad Karp. 2008. Wedge: Splitting applications into reduced-privilege compartments. In *Proceedings of the 5th USENIX Symposium on Networked Systems Design and Implementation (NSDI'08)*. USENIX Association, Berkeley, CA, 309–322. Retrieved from <http://dl.acm.org/citation.cfm?id=1387589.1387611>.
- [10] Daniel Bittman, Peter Alvaro, Darrell D. E. Long, and Ethan L. Miller. 2019. A tale of two abstractions: The case for object space. In *Proceedings of USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage'19)*.
- [11] Daniel Bittman, Peter Alvaro, and Ethan L. Miller. 2019. A persistent problem: Managing pointers in NVM. In *Proceedings of the 10th Workshop on Programming Languages and Operating Systems (PLOS'19)*. 30–37.
- [12] Adrian M. Caulfield, Arup De, Joel Coburn, Todor Mollov, Rajesh Gupta, and Steven Swanson. 2010. Moneta: A high-performance storage array architecture for next-generation, non-volatile memories. In *Proceedings of the 43rd IEEE/ACM International Symposium on Microarchitecture (MICRO'10)*. 385–395. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/caulfield-micro10.pdf>.
- [13] Jeffrey S. Chase, Henry M. Levy, Michael J. Feeley, and Edward D. Lazowska. 1994. Sharing and protection in a single-address-space operating system. *ACM Trans. Comput. Syst.* 12, 4 (Nov. 1994), 271–307. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/chase-tocs94.pdf>.
- [14] Guoyang Chen, Lei Zhang, Richa Budhiraja, Xipeng Shen, and Youfeng Wu. 2017. Efficient support of position independence on non-volatile memory. In *Proceedings of the 50th IEEE/ACM International Symposium on Microarchitecture (MICRO'17)*. ACM, New York, NY, 191–203. Retrieved from <http://doi.acm.org/10.1145/3123939.3124543>.
- [15] Howard Chu and Symas. [n.d.]. Lightning Memory-Mapped Database (part of the OpenLDAP project). Retrieved from <https://symas.com/lmdb/>.
- [16] Joel Coburn, Adrian M. Caulfield, Ameen Akel, Laura M. Grupp, Rajesh K. Gupta, Ranjit Jhala, and Steven Swanson. 2011. NV-Heaps: Making persistent objects fast and safe with next-generation, non-volatile memories. In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'11)*. 105–118. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/coburn-asplos11.pdf>.
- [17] Jeremy Condit, Edmund B. Nightingale, Christopher Frost, Engin Ipek, Benjamin Lee, Doug Burger, and Derrick Coetzee. 2009. Better I/O through byte-addressable, persistent memory. In *Proceedings of the 22nd ACM Symposium*

- on *Operating Systems Principles (SOSP'09)*. 133–146. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/condit-sosp09.pdf>.
- [18] Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. 2010. Benchmarking cloud serving systems with YCSB. In *Proceedings of the 1st ACM Symposium on Cloud Computing (SoCC'10)*. ACM, New York, NY, 143–154. DOI: <https://doi.org/10.1145/1807128.1807152>.
- [19] Fernando J. Corbató and Victor A. Vyssotsky. 1965. Introduction and overview of the Multics system. In *Proceedings of the November 30 – December 1, 1965, Fall Joint Computer Conference, Part I*. ACM, 185–196. Retrieved from <http://dl.acm.org/citation.cfm?id=1463912>.
- [20] Robert C. Daley and Jack B. Dennis. 1968. Virtual memory, processes, and sharing in MULTICS. *Commun. ACM* 11, 5 (May 1968), 306–312. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/daley-cacm68.pdf>.
- [21] Partha Dasgupta, Richard J. LeBlanc, Jr., Mustaque Ahamad, and Umakishore Ramachandran. 1991. The clouds distributed operating system. *IEEE Comput.* (Nov. 1991). Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/dasgupta-computer91.pdf>.
- [22] Alan Dearle, Rex di Bona, James Farrow, Frans Henskens, Anders Lindström, John Rosenberg, and Francis Vaughan. 1994. Grasshopper: An orthogonally persistent operating system. *Comput. Syst.* 7, 3 (June 1994), 289–312. Retrieved from <http://dl.acm.org/citation.cfm?id=198008.198009>.
- [23] Biplob Debnath, Sudipta Sengupta, and Jin Li. 2010. FlashStore: High throughput persistent key-value store. In *Proceedings of the 36th Conference on Very Large Databases (VLDB'10)*. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/debnath-vldb10.pdf>.
- [24] Xiangyu Dong, Cong Xu, Norm Jouppi, and Yuan Xie. 2014. *Emerging Memory Technologies: Design, Architecture, and Applications*. Springer, 15–50.
- [25] Subramanya R. Dulloor, Sanjay Kumar, Anil Keshavamurthy, Philip Lantz, Dheeraj Reddy, Rajesh Sankaran, and Jeff Jackson. 2014. System software for persistent memory. In *Proceedings of the 9th European Conference on Computer Systems (EuroSys'14)*. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/dulloor-eurosys14.pdf>.
- [26] Izzat El Hajj, Alexander Merritt, Gerd Zellweger, Dejan Milojicic, Reto Achermann, Paolo Faraboschi, Wen-mei Hwu, Timothy Roscoe, and Karsten Schwan. 2016. SpaceJMP: Programming with multiple virtual address spaces. In *Proceedings of the 21st International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'16)*. ACM, New York, NY, 353–368. DOI: <https://doi.org/10.1145/2872362.2872366>.
- [27] Dawson R. Engler, Sandeep K. Gupta, and M. Frans Kaashoek. 1995. AVM: Application-level virtual memory. In *Proceedings of the 5th Workshop on Hot Topics in Operating Systems (HotOS'95)*. IEEE, 72–77.
- [28] Dawson R. Engler and M. Frans Kaashoek. 1995. Exterminate all operating system abstractions. In *Proceedings of the 5th Workshop on Hot Topics in Operating Systems (HotOS'95)*. IEEE, 78–83.
- [29] Dawson R. Engler, M. Frans Kaashoek, and James O'Toole, Jr. 1995. Exokernel: An operating system architecture for application-level resource management. In *Proceedings of the 15th ACM Symposium on Operating Systems Principles (SOSP'95)*. 251–266. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/engler-sosp95.pdf>.
- [30] Hewlett Packard Enterprise. 2018. YCSB-C. Retrieved from <https://github.com/HewlettPackard/meadowlark/tree/master/extra/YCSB-C>.
- [31] Paolo Faraboschi, Kimberly Keeton, Tim Marsland, and Dejan Milojicic. 2015. Beyond processor-centric operating systems. In *Proceedings of the 15th Workshop on Hot Topics in Operating Systems (HotOS'15)*. USENIX Association. Retrieved from <https://www.usenix.org/conference/hotos15/workshop-program/presentation/faraboschi>.
- [32] David K. Gifford, Pierre Jouvelot, Mark A. Sheldon, and James W. O'Toole, Jr. 1991. Semantic file systems. In *Proceedings of the 13th ACM Symposium on Operating Systems Principles (SOSP'91)*. ACM, 16–25. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/gifford-sosp91.pdf>.
- [33] Burra Gopal and Udi Manber. 1999. Integrating content-based access mechanisms with hierarchical file systems. In *Proceedings of the 3rd Symposium on Operating Systems Design and Implementation (OSDI'99)*. 265–278. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/gopal-osdi99.pdf>.
- [34] Jorge Guerra, Leonardo Mármol, Daniel Campello, Carlos Crespo, Raju Rangaswami, and Jinpeng Wei. 2012. Software persistent memory. In *Proceedings of the USENIX Annual Technical Conference*. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/guerra-atc12.pdf>.
- [35] Gernot Heiser and Kevin Elphinstone. 2016. L4 microkernels: The lessons from 20 years of research and deployment. *ACM Trans. Comput. Syst.* 34, 1 (April 2016). DOI: <https://doi.org/10.1145/2893177>.
- [36] Gernot Heiser, Kevin Elphinstone, Stephen Russell, and Jerry Vochtelloo. 1993. *Mungi: A Distributed Single Address-Space Operating System*. Technical Report 9314. School of Computer Science and Engineering, University of New South Wales. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/heiser-scse9314.pdf>.
- [37] Benjamin Hindman, Andy Konwinski, Matei Zaharia, Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, and Ion Stoica. 2011. Mesos: A platform for fine-grained resource sharing in the data center. In *Proceedings of the*

- 8th USENIX Conference on Networked Systems Design and Implementation (NSDI'11). USENIX, Berkeley, CA, 295–308. Retrieved from <http://dl.acm.org/citation.cfm?id=1972457.1972488>.
- [38] Antony L. Hosking and J. Eliot B. Moss. 1993. Object fault handling for persistent programming languages: A performance evaluation. In *Proceedings of the 8th Conference on Object-oriented Programming Systems, Languages, and Applications (OOPSLA'93)*. ACM, New York, NY, 288–303. DOI : <https://doi.org/10.1145/165854.165907>.
- [39] Qingda Hu, Jinglei Ren, Anirudh Badam, and Thomas Moscibrod. 2017. Log-structured non-volatile main memory. In *Proceedings of the USENIX Annual Technical Conference*. 703–717. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/hu-atc17.pdf>.
- [40] Joseph Izraelevitz, Jian Yang, Lu Zhang, Juno Kim, Xiao Liu, Amirsaman Memaripour, Yun Joon Soh, Zixuan Wang, Yi Xu, Subramanya R. Dulloor, Jishen Zhao, and Steven Swanson. 2019. Basic performance measurements of the Intel Optane DC persistent memory module. *arXiv abs/1903.05714* (2019).
- [41] Eric Jul, Henry Levy, Norman Hutchinson, and Andrew Black. 1988. Fine-grained mobility in the Emerald system. *ACM Trans. Comput. Syst.* 6, 1 (Feb. 1988), 109–133. DOI : <https://doi.org/10.1145/35037.42182>.
- [42] Eric Jul and Bjarne Steensgaard. 1991. Implementation of distributed objects in Emerald. In *Proceedings of the International Workshop on Object Orientation in Operating Systems*. IEEE, 130–132.
- [43] M. Frans Kaashoek, Dawson R. Engler, Gregory R. Ganger, Hector M. Briceño, Russell Hunt, David Mazières, Thomas Pinckney, Robert Grimm, John Jannotti, and Kenneth Mackenzie. 1997. Application performance and flexibility on exokernel systems. In *Proceedings of the 16th ACM Symposium on Operating Systems Principles (SOSP'97)*. ACM, New York, NY, 52–65. DOI : <https://doi.org/10.1145/268998.266644>.
- [44] Orran Krieger, Marc Auslander, Bryan Rosenburg, Robert W. Wisniewski, Jimi Xenidis, Dilma Da Silva, Michal Ostrowski, Jonathan Appavoo, Maria Butrico, Mark Mergen, Amos Waterland, and Volkmar Uhlig. 2006. K42: Building a complete operating system. In *Proceedings of the 1st ACM SIGOPS/EuroSys European Conference on Computer Systems 2006 (EuroSys'06)*. ACM, New York, NY, 133–145. DOI : <https://doi.org/10.1145/1217935.1217949>.
- [45] Dokeun Lee and Youjip Won. 2013. Bootless boot: Reducing device boot latency with byte addressable NVRAM. In *Proceedings of the International Conference on High Performance Computing*. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/lee-hpcc13.pdf>.
- [46] James Litton, Anjo Vahldiek-Oberwagner, Eslam Elnikety, Deepak Garg, Bobby Bhattacharjee, and Peter Druschel. 2016. Light-weight contexts: An OS abstraction for safety and performance. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI'16)*. USENIX Association, 49–64. Retrieved from <https://www.usenix.org/conference/osdi16/technical-sessions/presentation/litton>.
- [47] Youyou Lu, Jiwu Shu, and Long Sun. 2016. Blurred persistence: Efficient transactions in persistent memory. *ACM Trans. Stor.* 12, 1 (Jan. 2016). Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/lu-tos16.pdf>.
- [48] Youyou Lu, Jiwu Shu, Long Sun, and Onur Mutlu. 2014. Loose-ordering consistency for persistent memory. In *Proceedings of the 32nd IEEE International Conference on Computer Design (ICCD'14)*. IEEE, 216–223.
- [49] Virendra J. Marathe, Margo Seltzer, Steve Byan, and Tim Harris. 2017. Persistent memcached: Bringing legacy code to byte-addressable persistent memory. In *Proceedings of the 9th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage'17)*. USENIX Association. Retrieved from <https://www.usenix.org/conference/hotstorage17/program/presentation/marathe>.
- [50] Pankaj Mehra and Samuel Fineberg. 2004. Fast and flexible persistence: The magic potion for fault-tolerance, scalability and performance in online data stores. In *Proceedings of the 18th International Parallel and Distributed Processing Symposium (IPDPS'04)*. DOI : <https://doi.org/10.1109/IPDPS.2004.1303232>.
- [51] Justin Meza, Yixin Luo, Samira Khan, Jishen Zhao, Yuan Xie, and Onur Mutlu. 2013. A case for efficient hardware/software cooperative management of storage and memory. In *Proceedings of the 5th Workshop on Energy-Efficient Design (WEED'13)*. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/meza-weed13.pdf>.
- [52] Dushyanth Narayanan and Orion Hodson. 2012. Whole-system persistence. In *Proceedings of the 17th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'12)*. 401–500. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/narayanan-asplos12.pdf>.
- [53] Yuanjiang Ni, Jishen Zhao, Daniel Bittman, and Ethan Miller. 2018. Reducing NVM writes with optimized shadow paging. In *Proceedings of the 10th Workshop on Hot Topics in Storage and File Systems (HotStorage'18)*. Retrieved from <http://www.ssrc.ucsc.edu/ni-hotstorage18.pdf>.
- [54] Yuanjiang Ni, Jishen Zhao, Heiner Litz, Daniel Bittman, and Ethan L. Miller. 2019. SSP: Eliminating redundant writes in failure-atomic NVRAMs via shadow sub-paging. In *Proceedings of the 52nd IEEE/ACM International Symposium on Microarchitecture*.
- [55] Matheus Ogleari, Ethan L. Miller, and Jishen Zhao. 2018. Steal but no force: Efficient hardware-driven undo+redo logging for persistent memory systems. In *Proceedings of the 24th International Symposium on High-performance Computer Architecture (HPCA'18)*. Retrieved from <http://www.ssrc.ucsc.edu/ogleari-hpca18.pdf>.

- [56] Yoann Padioleau and Olivier Ridoux. 2003. A logic file system. In *Proceedings of the USENIX Annual Technical Conference*. 99–112. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/padioleau-usenix03.pdf>.
- [57] Aleatha Parker-Wood, Darrell D. E. Long, Ethan L. Miller, Philippe Rigaux, and Andy Isaacson. 2014. A file by any other name: Managing file names with metadata. In *Proceedings of the 7th International Systems and Storage Conference (SYSTOR'14)*. Retrieved from <http://www.ssrc.ucsc.edu/Papers/parkerwood-systor14.pdf>.
- [58] Timothy Roscoe. 1994. Linkage in the Nemesis single address space operating system. *ACM SIGOPS Oper. Syst. Rev.* 28, 4 (Oct. 1994), 48–55. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/roscoe-osr94.pdf>.
- [59] Andy Rudoff et al. 2017. Persistent Memory Programming Library. Retrieved from <http://pmem.io/nvml/>.
- [60] Yizhou Shan, Shin-Yeh Tsai, and Yiyang Zhang. 2017. Distributed shared persistent memory. In *Proceedings of the Symposium on Cloud Computing (SoCC'17)*. Association for Computing Machinery, New York, NY, 323–337. DOI : <https://doi.org/10.1145/3127479.3128610>.
- [61] Jonathan S. Shapiro and Jonathan Adams. 2002. Design evolution of the EROS single-level store. In *Proceedings of the USENIX Annual Technical Conference*. USENIX, 59–72. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/shapiro-usenix02.pdf>.
- [62] Jonathan S. Shapiro, Jonathan M. Smith, and David J. Farber. 1999. EROS: A fast capability system. In *Proceedings of the 17th ACM Symposium on Operating Systems Principles (SOSP'99)*. ACM, New York, NY, 170–185. DOI : <https://doi.org/10.1145/319151.319163>.
- [63] Eugene Shekita and Michael Zwilling. 1990. *Cricket: A Mapped, Persistent Object Store*. Technical Report 956. University of Wisconsin. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/shekita-uw-tr956.pdf>.
- [64] Alan Skousen and Donald Miller. 1999. Using a single address space operating system for distributed computing and high performance. In *Proceedings of the 18th IEEE International Performance, Computing and Communications Conference (IPCCC'99)*. 8–14.
- [65] Hung-Wei Tseng, Qianchen Zhao, Yuxiao Zhou, Mark Gahagan, and Steven Swanson. 2016. Morpheus: Creating application objects efficiently for heterogenous computing. In *Proceedings of the ACM/IEEE 43rd International Symposium on Computer Architecture*.
- [66] Haris Volos, Andres Jaan Tack, and Michael M. Swift. 2011. Mnemosyne: Lightweight persistent memory. In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'11)*. Retrieved from <http://www.ssrc.ucsc.edu/PaperArchive/volos-asplos11.pdf>.
- [67] Tiancong Wang, Sakthikumaran Sambasivam, Yan Solihin, and James Tuck. 2017. Hardware supported persistent object address translation. In *Proceedings of the 50th IEEE/ACM International Symposium on Microarchitecture (MICRO'17)*. ACM, New York, NY, 800–812. DOI : <https://doi.org/10.1145/3123939.3123981>.
- [68] William Wulf, Ellis Cohen, William Corwin, Anita Jones, Roy Levin, C. Pierson, and Fred Pollack. 1974. HYDRA: The kernel of a multiprocessor operating system. *Commun. ACM* 17, 6 (June 1974), 337–345. DOI : <https://doi.org/10.1145/355616.364017>.
- [69] Jian Xu and Steven Swanson. 2016. NOVA: A log-structured file system for hybrid volatile/non-volatile main memories. In *Proceedings of the 14th Usenix Conference on File and Storage Technologies (FAST'16)*. USENIX Association, Berkeley, CA, 323–338. Retrieved from <http://dl.acm.org/citation.cfm?id=2930583.2930608>.

Received November 2020; revised March 2021; accepted March 2021