

Low-Latency Synchronous I/O For OpenZFS Using Persistent Memory

Masterarbeit
von

Christian Schwarz

an der Fakultät für Informatik

Erstgutachter:	Prof. Dr. Frank Bellosa
Zweitgutachter:	Prof. Dr. Wolfgang Karl
Betreuender Mitarbeiter:	Lukas Werling, M.Sc.

Bearbeitungszeit: TODO – TODO

Ich versichere wahrheitsgemäß, die Arbeit selbstständig verfasst, alle benutzten Hilfsmittel vollständig und genau angegeben und alles kenntlich gemacht zu haben, was aus Arbeiten anderer unverändert oder mit Abänderungen entnommen wurde sowie die Satzung des KIT zur Sicherung guter wissenschaftlicher Praxis in der jeweils gültigen Fassung beachtet zu haben.

Karlsruhe, den 19. Mai 2021

Abstract

We propose to explore persistent memory (PMEM) as the storage medium for the ZFS Intent Log (ZIL). Specifically, we propose a new type of ZIL called ZIL-PMEM that bypasses existing block-device oriented abstractions to take advantage of the low latency provided by PMEM. ZIL-PMEM will maintain the same consistency guarantees to applications as the current ZIL implementation, maintain data integrity of log contents and correctly handle data corruption.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special contents, but the length of words should match the language.

Contents

Abstract	v
Contents	1
1 Introduction	5
2 Literature Review & Background	9
2.1 Literature Review	9
2.1.1 PMEM Filesystems	9
2.1.2 Cross-Media Systems	12
2.1.3 Journaling & Write-Ahead Logs Adapted PMEM	16
2.1.4 Testing Filesystem Crash Consistency	17
2.1.5 PMEM-specific Crash-Consistency Checkers	19
2.1.6 Fault Injection	20
2.2 Persistent Memory in Linux	22
2.3 OpenZFS Primer	22
2.3.1 The ZIL API	22
2.3.2 Replay	24
2.3.3 Log Record Types	24
3 ZIL-LWB on PMEM	27
3.1 Benchmark Setup	27
3.2 Results	29
3.3 Analysis	31
3.3.1 OpenZFS Background: How ZIL-LWB Works	31
4 Design Overview	39
4.1 Project Goals & Scope	39
4.1.1 Requirements	39
4.1.2 Out Of Scope For The Thesis	41
4.1.3 Limitations	41

4.2	Design Overview	42
5	The PRB/HDL Data Structure	45
5.1	OpenZFS Background	45
5.2	Requirements	48
5.3	Approach	49
5.4	HDL: Log Structure	50
5.5	HDL Replay: Basic Approach	52
5.6	HDL Replay: Dependency Tracking	53
5.7	HDL: Model For Data Corruption	58
5.8	HDL: Replay Crash-Consistency	59
5.9	PRB: DRAM Data Structure	61
5.10	PRB: PMEM Data Structure	65
5.11	PRB: Chunk Traversal	66
5.12	PRB: Garbage Collection	69
5.13	The Write Path	69
5.13.1	Commit Slots	71
5.13.2	HDL-Scoped Metadata	73
5.13.3	Crash-Consistent Insert	74
5.14	API Design	77
5.14.1	PRB Setup	77
5.14.2	Replay	79
5.14.3	Writing Entries	80
6	Integration into ZFS	83
6.1	ZIL Kinds	83
6.1.1	On-Disk State	84
6.1.2	Runtime State	85
6.1.3	Changing ZIL Kinds	89
6.1.4	ZIL-LWB Suspend & Resume	89
6.1.5	ZIL Traversal & ZDB	90
6.1.6	ZIL-LWB-Specific Callbacks	91
6.1.7	Reflection & Alternatives	92
6.2	PMEM-aware SPA & VDEV layer	94
6.3	The ZIL-PMEM ZIL Kind	95
6.3.1	Activation	95
6.3.2	PRB Construction	95
6.3.3	Dataset & HDL Lifecycle Synchronization	96
6.3.4	ZILOG_PMEM_T	97
6.3.5	ITXG Bypass For ZVOL	100

7	Evaluation	103
7.1	Usability & Architecture	103
7.2	Correctness	104
7.2.1	PRB	104
7.2.2	ZIL-PMEM	105
7.2.3	ZIL-PMEM	108
7.3	Performance	110
7.3.1	Setup & Methodology	110
7.3.2	4k Synchronous Random Write Workload	110
7.3.3	Application Benchmarks	113
7.3.4	ZIL-PMEM vs. ZIL-LWB	118
7.3.5	ZIL-PMEM vs. XFS and Ext4 on PMEM	119
7.3.6	ZIL-PMEM vs. XFS and Ext4 on Dm-writocache	121
7.3.7	Impact on ZVOL Performance & ITXG Bypass	123
7.3.8	CPU-Efficient Handling Of PMEM Bandwidth Limits	125
7.4	Discussion	129
8	Summary	131
8.1	Future Work	131
	Appendix	133
	Bibliography	135

Chapter 1

Introduction

The task of a filesystem is to provide non-volatile storage to applications in the form of the *file* abstraction. Applications modify files through system calls such as `write()` which generally does not provide any durability guarantees. Instead, the system call modifies a buffer in DRAM such as a page in the Linux page cache and returns to userspace. Synchronization of the dirty in-DRAM data to persistent storage is thus deferred to a — generally implementation-defined — point in the future.

However, many applications have more specific durability requirements. For example, an accounting system that processes a purchase needs to ensure that the updated account balance is persisted to non-volatile storage before clearing the transaction. Otherwise, a system crash and recovery after clearing could result in the pre-purchase balance being restored, enabling double-spending by the account holder. These **synchronous I/O** semantics must be requested through APIs such as `fsync()` which “assure that after a system crash [...] all data up to the time of the `fsync()` call is recorded on the disk.” [[posix_fsync_opengroup](#)].

The **Zettabyte File System (ZFS)** is a combined volume manager and filesystem. It pools many block devices into a single storage pool (*zpool*) which can hold thousands of sparsely allocated filesystems. The ZFS on-disk format is a merkle tree that is rooted in the *uberblock* which is ZFS’s equivalent of a superblock. ZFS on-disk state is always consistent and moves forward in so-called *transaction groups* (txg), using copy-on-write to apply updates. Whenever a new version of the on-disk state needs to be synced to disk, ZFS traverses its logical structure bottom up and builds a new merkle tree. The updated parts of the tree are stored in newly allocated disk blocks while unmodified parts are re-use the existing block written in a prior txg. Once all updates have been written out, the new

uberblock is written, thereby atomically moving the on-disk format to its new state. This procedure is called *txg sync* and is triggered periodically (default: every 5 s) or if the amount of dirty data exceeds a configurable threshold.

Synchronous I/O semantics cannot be reasonably implemented through *txg sync* due to the write amplification and CPU overhead inherent to the *txg sync* procedure. Instead ZFS maintains the **ZFS Intent Log (ZIL)** which is a per-filesystem write-ahead log. Unlike systems such as Linux's *journaling block device 2* (JBD2), the ZIL is a logical log: the ZIL's records describe the *logical* changes that need to be applied in order to achieve the state that was reported committed to userspace. On disk, the log records are written into a chain of *log-write blocks* (LWBs), each containing many records. The LWB chain is rooted in the *ZIL header* within the filesystem's representation within the merkle tree. New LWBs are appended to the chain independently of *txg sync*.

from?

By default, LWBs are allocated on the zpool's main storage devices. Consequently, the lower bound for synchronous I/O latency in ZFS is the time required to write the LWBs that contains the synchronous I/O operation's ZIL records. For the case where this latency is insufficient, ZFS provides the ability to add a *separate log device* (SLOG) to the zpool. The SLOG is typically a single (or mirrored) block device that provides lower latency than the main pool's devices. A typical configuration today is to add a fast NVMe drive to an HDD-based pool. Adding a fast SLOG accelerates LWB writes because LWBs are preferentially allocated from the SLOG. Note that SLOGs only need very limited capacity since LWBs are generally obsolete after three txgs.

Persistent Memory (PMEM) is an emerging storage technology that provides low-latency memory-mapped byte-addressable persistent storage. The Linux kernel can expose PMEM as a pseudo block device whose sectors map directly to the PMEM space. Thereby existing block devices consumers can benefit from PMEM's low latency without modification. Block device consumers that wish to bypass block device emulation use the kernel-internal *DAX* API which translates sector numbers to kernel virtual addresses, giving software **direct access** to PMEM.

product name

The motivation for this thesis is to accelerate synchronous I/O in ZFS by using PMEM as a ZFS SLOG device. A single DIMM of the current Intel Optane PMEM product line can sustain 530k random 4k write IOPS which corresponds to a write latency of 1.88 μ s. However, when configuring the Linux PMEM block device as a SLOG in ZFS, a single thread only achieves 12k random 4k synchronous write IOPS ($\sim 83 \mu$ s), scaling up to 100k IOPS at 16 threads (8 cores, 2xSMT) at doubled latency. In contrast, the same workload applied directly to the PMEM block device is able to achieve 500k IOPS with two threads ($\sim 4 \mu$ s per thread).

And Linux 5.9's ext4 on the PMEM block device achieves 100k random 4k write IOPS (\sim) with a single thread and scales approximately linearly to up to 466k IOPS at four threads.

We perform a more detailed analysis to determine ZFS's severe latency overhead in the 4k sync write benchmark. We find that the 80% of the wall clock time for the average write operation is spent in the ZIL code. We identify severe overheads caused by both the physical data structure that represents the ZIL as well as the mechanism to persist it.

Based on these observations, we propose **ZIL-PMEM**, a new ZIL design that exclusively targets PMEM. It coexists with the existing ZIL which we refer to as *ZIL-LWB* in the remainder of this document. ZIL-PMEM achieves 140k random 4k sync write IOPS with a single thread (~ 7.1 us) and scales up to 400k IOPS at eight threads before it becomes CPU-bound. This corresponds to a speedup of 6.8x (or 4x, respectively) over ZIL-LWB. Our implementation is extensively unit-tested and passes the ZFS test suite's SLOG integration tests.

The remainder of this thesis is structured as follows: in Chapter 2 we review prior work in the field of persistent memory related to filesystems and provide background knowledge on the integration of PMEM in the Linux kernel as well as a detailed introduction to the components of ZFS that are relevant for this thesis. Chapter 3 describes our analysis of ZIL-LWB performance with a PMEM SLOG, motivating a new PMEM-specific ZIL implementation. Our solution, ZIL-PMEM, is then presented Chapters 4, 5, and 6. Chapter 4 defines the project goals and provides a high-level overview of the overall design. Our main contribution, the *PRB/HDL* the data structure, is then presented in Chapter 5. Chapter 6 then presents how we integrate integrate PRB/HDL into ZFS. We evaluate our implementation in Chapter 7 and conclude with a summary of our work in Chapter 8.

review numbers in this paragraph

review numbers

Chapter 2

Literature Review & Background

In this chapter we present prior work in the field of persistent memory and its application in various storage systems developed in research and industry. The last two sections on PMEM in Linux and ZFS provide the technical background knowledge that is necessary to understand the performance analysis of ZIL-LWB in Chapter 3 and the design of ZIL-PMEM in Chapter ??.

2.1 Literature Review

We have surveyed publications in the area of persistent memory storage systems, filesystem guarantees & crash-consistency models, PMEM-specific crash-consistency checkers and general methods to determine filesystem robustness in the presence of hardware failures.

2.1.1 PMEM Filesystems

In this subsection we present research filesystems that were explicitly designed for persistent memory. ZIL-PMEM integrates into ZFS, a production filesystem that was not designed for persistent memory. Hence we focus on techniques for crash consistency and data integrity that might be applicable to our work.

In-Kernel PMEM Filesystems

The initial wave of publications around the use of PMEM in filesystems produced a set of systems that were implemented completely in the kernel.

BPFS [Con+09] is one of the earliest filesystems expressly designed for PMEM. The filesystem layout in PMEM is inspired by WAFL ([HLM94]) and resembles a

we don't yet compare these systems with ZIL-PMEM. should we? Maybe in a compact section after we presented the design?

tree of multi-level indirect pages that eventually point to data pages. BPFS's key contribution is the use of fine-grained atomic updates in lieu of journaling for crash consistency. For example, updates to small metadata such as *mtime* can be made using atomic operations. For larger modifications, the authors introduce *short-circuit shadow-paging*, a technique where updates are prepared in a copy of the page. The updated page is then made visible through an update of the pointer in its parent indirect page. The difference to regular copy-on-write is that, as soon as the update to an indirect page can be done through an atomic in-place operation, the atomic operation is used. Updates thereby do not necessarily propagate up to the root of the tree.

PMFS [Dul+14] is another research filesystem that targets persistent memory. The authors make frequent comparisons to BPFS. The main differentiator from BPFS with regards to consistency is PMFS's use of undo-logging for metadata updates and copy-on-write for data consistency in addition to hardware-provided atomic in-place updates. The evaluation shows that their approach for metadata has between 24x and 38x lower overhead compared to BPFS (unit: number of bytes copied). PMFS also introduces an efficient protection mechanism against accidental *scribbles*. Scribbles are bugs in the system that accidentally overwrite PMEM, e.g., due to incorrect address calculation or out-of-bounds access in the kernel. By default, PMFS maps all PMEM read-only, thereby preventing accidental corruption of PMEM from code outside of PMFS. When PMFS needs to modify PMEM, it temporarily disables interrupts and clears the processor's CR0.WP, thereby allowing writes to read-only pages in kernel mode. [Dul+14; SDM413]

NOVA [XS16] is the most mature research PMEM filesystem. NOVA uses per-inode logs for operations scoped to a single inode (e.g. write syscalls) and per-CPU journals for operations that affect multiple inodes. The intended result is high scalability with regard to core count. The per-inode log data structure is a linked list with a head and tail pointer in the inode. NOVA leverages 8-byte atomic operations to update these pointers after it has written log entries. While not explicitly called so by the authors it is our impression that the log is a logical redo log except for writes, which — judging from the text — are always logged at page granularity. The authors explain the recovery procedure and measure its performance but do not address correctness in the evaluation.

NOVA-Fortis [Xu+17] is a version of NOVA that introduces snapshots and hardening against data corruption. Whereas snapshots are not relevant for this thesis because ZFS already provides this feature, the data corruption countermeasures are representative of the state of the art:

- Handling of machine check exceptions (MCE) in case the hardware detects bit errors. This is done by using `memcpy_mcsafe()` for all PMEM access.

- Detection of metadata corruption through CRC32 checksums.
- Redundancy through metadata replication. For metadata recovery, NOVA-Fortis compares checksums of the primary and replica and restores the variant with the matching checksum.
- Protection against localized unrecoverable data loss through RAID4-like parity. This feature only works while the file is not DAX-mapped.
- Protection against Scribbles using `CR0.WP` as described in the paragraph on PMFS.

The authors use a custom fault injection tool to corrupt data structures in a targeted manner and test NOVA Fortis’s recovery capabilities.

Hybrid PMEM filesystems

A recurring pattern in PMEM filesystem design is to split responsibilities between kernel and userspace component in order to eliminate system call overhead.

Aerie [Vol+14] is a user-space filesystem based on the premise that “SCM [Storage Class Memory] no longer requires the OS kernel [...]. Applications link to a file-system library that provides local access to data and communicates with a [user-space] service for coordination. The OS kernel provides only coarse grained allocation and protection, and most functionality is distributed to client programs. For read-only workloads, applications can access data and metadata through memory with calls to the file-system service only for coarse-grained synchronization. When writing data, applications can modify data directly but must contact the file-system service to update metadata.” The system uses a redo log maintained in each client program which is shipped to the filesystem service periodically or when a global lock is released. It is our understanding that only log entries shipped to and validated by the filesystem service will be replayed. The authors state that log entries can be lost if a client crashes before the log entries are shipped. The evaluation does not address crash consistency or recovery at all.

Strata [Kwo+17] is a cross-media filesystem with both kernel and user-space components. Since its distinguishing feature is the intelligent migration of data between different storage media we discuss it in the Section 2.1.2 on cross-media storage systems.

SplitFS [Kad+19] is a research filesystem that proposes a “split of responsibilities between a user-space library file system and an existing kernel PM file system. The user-space library file system handles data operations by intercepting POSIX calls, memory-mapping the underlying file, and serving the read and

overwrites using processor loads and stores. Metadata operations are handled by the kernel PM file system (ext4 DAX)”. SplitFS uses a redo log with idempotent entries that is written from userspace. As a performance optimization the authors use checksums and aligned start addresses to find valid log entries instead of an explicit linked list with persistent pointers. The SplitFS evaluation of correctness is limited to a comparison of user-observable filesystem state between SplitFS and ext4 in DAX mode. Recovery is evaluated only through the lens of recovery time (performance), not correctness.

EvFS [YCH19] is a “user-level POSIX file system that directly manages NVM in user applications. EvFS minimizes the latency by building a user-level storage stack and introducing asynchronous processing of complex file I/O with page cache and direct I/O. [...] EvFS leads to a 700-ns latency for 64-byte non-blocking file writes and reduces the latency for 4-Kbyte blocking file I/O by 20 us compared to a kernel file system [EXT4] with journaling disabled.” In contrast to Aerie, EvFS does not require a coordinating user-space service. Crash consistency and recovery is not addressed: “EvFS is not a production-ready file system because it neither provides all the POSIX APIs or crash-safe properties”.

2.1.2 Cross-Media Systems

Cross-media storage systems combine the advantages of multiple storage devices from different levels of the storage hierarchy. Historically, these kinds of systems strive to exploit hard disk for high capacity at low cost and more expensive flash storage for low latency random access IO. With persistent memory, a new class of storage has become available whose role in cross-media systems is still to be determined. In the context of this thesis we find it most useful to compare the overall system architecture.

ZFS: Allocation Classes [] Whereas ZFS was initially designed for a large pool of hard disks, it has gained several cross-media features over its lifetime. When configuring a zpool, the administrator assigns the block device to an *allocation class*. The following allocation classes exist: *normal* (main pool), *log* (SLOG device), *aux* (L2ARC, a victim cache for ZFS’s ARC), *special* (small blocks), and *dedup* (deduplication table data). When a function in ZFS allocates a block in the pool, it must specify the desired allocation class. If the allocation succeeds, the allocated block is guaranteed to be located on device(s) within that class. The use case for allocation classes is to combine the advantages of different storage media in a single pool. In many setups devices in the *normal* class have high capacity and are grouped in storage-efficient redundancy configurations such as *raidz* or *draid*. A *log* or *aux* device can be added to a pool in order to accelerate latency-sensitive I/O such as ZIL writes. Many administrators configure lower

check that we have explained ARC already

redundancy for these devices — either to gain performance or to reduce costs — which, depending on business requirements, may be justified due to the short-lived nature of ZIL writes. In comparison to the other systems presented in this section, allocation classes are very inflexible: once an allocation is made and the data is written, the data stays in that place until it is freed or the device is removed from the pool. In particular, there is no automatic tiering that takes usage patterns into account. Further, L2ARC devices reduce space efficiency because the cached data still occupies space in the main pool. Similarly, SLOG devices are somewhat wasteful since they are exclusively used for ZIL logging and never read except during recovery. Systems such as Strata derives more value from its PMEM-based log. Finally, it should be noted that while Linux’s `/dev/pmem block` device can be used with allocation classes, ZFS’s ZIO pipeline is unable to exploit its write performance..

ref backwards

ref section on
ZIL-LWB perf

ZFS: ZIL Performance Improvements For Fast Media [Ope20] At the Open-ZFS 2020 Developer summit, Saji Nair of storage vendor Nutanix presented a ZIL prototype that handles fast block devices more efficiently. The prototype avoids unnecessary sequential ordering of I/O operations when writing ZIL LWBs. It also avoids using the ZIO pipeline due to context switching overheads, issuing block I/O directly from the application thread instead. The evaluation is limited to 4k sync write performance, claiming an up 4x improvement in IOPS with four threads. However, the source code for the prototype has not been published and the design is incomplete with regards to replay.

check that
LWBs have
been intro-
duced by now

Strata [Kwo+17] is a cross-media research filesystem. “Closest to the application, Strata’s user library synchronously logs process-private updates in NVM while reading from shared, read-optimized, kernel-maintained data and meta-data. [...] Client code uses the POSIX API, but Strata’s synchronous updates obviate the need for any sync-related system calls”. We classify Strata as a hybrid filesystem in Section 2.1.1 because it consists of both a userspace library and an in-kernel component. The log, written from userspace, is an idempotent logical redo log. The kernel component then *digests* the logs asynchronously, performing aggregation of the logged operations during digestion. Aggregation permits the kernel component to issue “sequential, aligned writes” to the slower storage

maybe this en-
tire section
should move
to the 'back-
ground' sec-
tion?

devices such as SSDs or HDDs. ZFS with and without ZIL-PMEM compares to Strata in the following ways:

- Both systems use a logical redo operation log instead of a block-level journaling mechanism.
- Both systems perform asynchronous write-back and thereby reap similar benefits from it (parallel batch processing, optimized allocation).
- Strata’s kernel component digests the logs written from user-space in order to write them back to other tiers. ZFS accumulates the write-back state in DRAM and never reads the log except for recovery. (It is unclear to us how often Strata digests the logs. ZFS performs write-back of dirty state after at most 5 seconds.)
- Strata seems to tune allocations for SSDs, e.g. allocating blocks in erasure block size to prevent write amplification. ZFS supports TRIM and supports variable block sizes up to 16 MiB, but there are no automatic optimizations that specifically target write amplification in SSDs.
- ZFS is in-kernel and requires no modifications to applications whereas Strata requires linking to or LD_PRELOADing a user-space library, which, tangentially, makes it incompatible with statically linked binaries.

Ziggurat [ZHS19] “Ziggurat exploits the benefits of NVMM through intelligent data placement during file writes and data migration. Ziggurat includes two placement predictors that analyze the file write sequences and predict whether the incoming writes are both large and stable, and whether updates to the file are likely to be synchronous. Ziggurat then steers the incoming writes to the most suitable tier based on the prediction: writes to synchronously-updated files go to the NVMM tier to minimize the synchronization overhead. Small, random writes also go to the NVMM tier to fully avoid random writes to disk. The remaining large sequential writes to asynchronously-updated files go to disk”. The authors compare Ziggurat to Strata as follows: “Strata is a multi-tiered user-space file system that exploits NVMM as the high-performance tier, and SSD/HDD as the lower tiers. It uses the byte-addressability of NVMM to coalesce logs and migrate them to lower tiers to minimize write amplification. File data can only be allocated in NVMM in Strata, and they can be migrated only from a faster tier to a slower one. The profiling granularity of Strata is a page, which increases the bookkeeping overhead and wastes the locality information of file accesses”. ZFS with and without ZIL-PMEM compares to Ziggurat as follows:

- Ziggurat actively migrates data into PMEM based on access pattern. ZFS has no provisions for data migration within the pool after block allocation.

The ZFS architecture such a feature is unlikely to be developed in the future (keyword: blockpointer rewrite).

- Ziggurat sends writes directly to the suitable tier based on prediction of future access patterns. If the access pattern is anticipated to be synchronous, Ziggurat choses the PMEM tier. The ZIL, and ZIL-PMEM specifically, only serves as a stop-gap between txg sync points of the main pool. Data is always written twice — once to the log and once to the main pool —, and never read from the log except during recovery.
- Both systems are fully in-kernel and require no modifications to user-space applications.
- Ziggurat builds on NOVA-Fortis and thus inherits the PMEM-specific data integrity and redundancy mechanisms provided for the PMEM storage tier. ZIL-PMEM data integrity measures are more limited but can be expanded in the future (see Chapter ??).
- The Ziggurat paper does not mention data integrity measures or redundancy mechanisms for the block device layers beneath PMEM. Possibly, existing Linux features such as Device Mapper could be used to compensate. In contrast, ZIL-PMEM benefits from ZFS’s strong data integrity and redundancy mechanisms once the logged data is txg synced.
- The Ziggurat design is “fast-first. It should use disks to expand the capacity of NVMM rather than using NVMM to improve the performance of disks as some previous systems have done”. ZIL-PMEM approaches persistent memory from the opposite direction, starting with a filesystem strongly focussed on block devices and only leveraging PMEM where appropriate.
- Ziggurat was not evaluated on actual persistent memory. The authors used memory on another NUMA node to simulate lower latencies. We evaluate ZIL-PMEM on commercially available PMEM hardware.

ensure txg sync has already been explained

dm-writocache [] is a new (2018) Linux Device Mapper target that implements a write-back caching layer for block devices. Given an *origin* device and a *cache* device, it exposes a new virtual block device with the same capacity as the origin device. Write-back is controlled by relative thresholds on the cache device space utilization (low and high watermark, default to 45% and 50%). Reads are not handled by dm-writocache and instead served from the Linux page cache. [] Dm-writocache is relevant to ZIL-PMEM because it is the closest functional analogue to the ZIL in the Linux Device Mapper stack. It has been designed with explicit support for persistent memory and its authors have presented it as a means to accelerate database workloads [Tad+19]. However, both the data flow and guarantees differ substantially between the two systems. First, ZIL-PMEM pools the PMEM space for all filesystems and ZVOLS in a zpool whereas dm-writocache is limited to a single block device and thus would require partition-

ing. Second, ZIL-PMEM persists writes as log entries to an append-only log structure and does not perform coalescing or delta-encoding. In contrast, *dm-writecache* is a block-level cache where an overwrite of the virtual block results in an over-write in the cache. Further, ZIL-PMEM protects data integrity through checksums whereas *dm-writecache* fully relies on hardware error correction and detection, reported via `memcpy_mcsafe()`. And after ZFS has synced out modifications to the main pool, they benefit from ZFS's strong data redundancy mechanisms such as *raidz*. In contrast, with *dm-writecache*, the origin block device driver must handle data redundancy if desired, e.g., through another Device Mapper target such as *dm-raid*. We compare the performance of the two implementations in Section 7.3.6 of our evaluation.

2.1.3 Journaling & Write-Ahead Logs Adapted PMEM

The following publications use persistent memory to accelerate file system journals and write-ahead logs.

Unioning of the Buffer Cache and Journaling Layers with Non-volatile Memory [LBN13] is an academic paper that presents a PMEM-aware buffer cache design which “subsumes the functionality of caching and [block-level] journaling”. From the buffer cache’s perspective, the on-disk blocks that make up a journal (e.g. one managed by JBD2) are indistinguishable from non-journal filesystem blocks. Thus, for a journal block A' that logs an update to a block A , both blocks A' and A will sit in the buffer cache. This waste of space can be reduced with PMEM by the author’s proposal. Instead of journaling on top of the block layer, one journals in the buffer cache itself by a) placing the buffer cache in PMEM and b) introducing a new buffer cache entry state “frozen” so that an entry can be “clean”, “dirty” or “frozen”. When modifying a block A , the filesystem no longer makes a journal entry but instead modifies the buffer cache entry directly. If the entry was “clean”, it is now “dirty”. If the entry was “frozen”, a copy is made and the modification goes to the copy, making it “dirty” as well. Committing a block to the journal is a simple state transition from “dirty” to “frozen”. On a crash + restart, “dirty” buffer cache entries are discarded but “frozen” entries remain. We find the approach exceptionally interesting and innovative and can imagine an application of the idea to the ZFS ARC. However, the system’s real-world performance is unclear (evaluation on DRAM). Also, we see non-trivial software engineering and maintenance problems with the approach. Finally, the paper does not address hardware error handling or data corruption concerns at all.

ext4 fast commits [] is a feature released in Linux 5.10. “The fast-commit journal [...] contains changes at the file level, resulting in a more compact format.

Information that can be recreated is left out, as described in the patch posting: For example, if a new extent is added to an inode, then corresponding updates to the inode table, the block bitmap, the group descriptor and the superblock can be derived based on just the extent information and the corresponding inode information. [...] Fast commits are an addition to — not a replacement of — the standard commit path; the two work together. If fast commits cannot handle an operation, the filesystem falls back to the standard commit path.” The cited article mentions ongoing work to use persistent memory for ext4 fast commits. The approach is inspired by per-inode journaling as proposed by Park and Shin in [PS17].

Disk-oriented database management systems often use write-ahead logs (WALs) to allow a transaction to commit before the modified pages are written back to stable storage. However, appending to a shared WAL file has historically been a latency and scalability bottleneck which lead to amortization techniques such as *commit groups* (aka *group commit*) and *pre-committed transactions* that batch the WAL entries of multiple committing transactions into a single physical WAL record.¹ With persistent memory, the latency bottleneck no longer lies in the raw I/O but rather the coordination overhead between multiple CPU cores and sockets, prompting new distributed logging designs that avoid a central point of contention. [Fan+11; Pel+13; Joh+10] ZIL-LWB employs *group commit* as well although the terminology is different. It batches log records issued by independent threads for the same filesystem into the currently *open log-write block* (LWB). Once the open LWB is full or a timeout has passed, the open LWB is written (*issued*) to disk. In contrast, ZIL-PMEM commits log records directly to PMEM and allows multiple cores to do so in parallel with minimal global coordination. Our evaluation shows that the design scales well on single-socket systems. NUMA and multi-socket systems, which are addressed in the cited publications from the database community, have been explicitly out of scope for this thesis.

2.1.4 Testing Filesystem Crash Consistency

In this section we survey techniques to determine and verify crash consistency guarantees of file systems. A formal and automatically verifiable model would have been extremely helpful to ensure that ZIL-PMEM maintains the same crash consistency guarantees as ZIL-LWB.

¹According to [DeW+84] “the notion of group commits appears to be past of the unwritten database folklore. The System-R implementors claim to have implemented it.”

ensure this term has been established by now

refer to section that explains the perf problems with ZIL-LWB

minimal global coordination = the commit slots + semaphore

meh, ref

ref

All File Systems Are Not Created Equal [Pil+14] contributes a survey of the atomicity and ordering guarantees of several popular Linux filesystems and provides a tool called ALICE to validate or derive the guarantees required by applications. ZFS with and without ZIL-PMEM could benefit from this survey as well. The survey could be used to characterize ZFS’s guarantees. ALICE could be used to determine whether ZFS’s guarantees are sufficient for applications and/or whether ZFS’s guarantees exceed the requirements of the majority of applications. However, since ZIL-PMEM shall maintain the same guarantees as ZIL-PMEM (see Section 4.1.1), such findings would only be relevant for future work.

Specifying and Checking File System Crash-Consistency Models [Bor+16]

The authors “present a formal framework for developing crash-consistency models, and a toolkit, called FERRITE, for validating those models against real file system implementations.” The system provides means to express expected filesystem behavior as a *litmus test*. A litmus test encodes expected behavior of a filesystem though a series of events (e.g. write to a file) and a final predicate expressed as a satisfiability problem. FERRITE can execute the litmus test against an axiomatic formal model of the filesystem to ensure that, iff the actual filesystem adheres to the formal model, the litmus test’s expectations hold. Notably the litmus test is executed symbolically and the validation predicates are checked for satisfiability by an SMT solver; this is exhaustive and not comparable to a unit or regression test. FERRITE can also execute litmus tests against the actual filesystem to test whether it adheres to the formal model. This test is non-exhaustive. It is based on executing the litmus test “many times” where for each execution all disk commands emitted during execution are recorded. All permutations and prefixes of these traces that are allowed under the semantics of the disk protocol are used to produce test disk images which are fed back to the filesystem under test for recovery. After recovery the litmus test’s predicate must hold against the concrete filesystem state after recovery. If it does not hold the given permutation of the trace is proof that the filesystem does not match its model (assuming the litmus test passes execution against the model), or that the filesystem’s assumptions about the disk do not match the FERRITE *disk model*. FERRITE appears to be a useful tool to build a model of ZFS’s undocumented crash-consistency guarantees (ZFS has not been evaluated by the authors). Such a model would be helpful to validate ZIL-PMEM’s goal to maintain the same semantics as ZIL-PMEM. However, this would require a FERRITE disk model for persistent memory.

Using Model Checking to Find Serious File System Errors [Yan+06]

The authors present an “implementation-level model checker” that removes the need to define a formal model for the filesystem. Instead, the model is inferred by running the OS with the filesystem and recording both syscalls emitted by the

application and disk operations emitted by the filesystem. These traces are subsequently fed to a process that produces disk images created from reorderings of disk operations. The disk images are fed to the filesystem’s fsck tool. Disk images that can be ‘repaired’ by fsck are then fed to the “recovery checker” which examines filesystem state and compares it to the expected state which (if we understand section 4.2 of the paper correctly) is derived from the recorded system calls. (The role of the “volatile file system” in this process is still unclear to us). The authors mention several shortcomings of their system:

- Lack of multi-threading support (this applies to FERRITE as well).
- The recovery checker produces a projection of the filesystem state (e.g. only names and content but no atime). Thus, the system can only check guarantees at the projection level.
- Restrictive assumptions such as “Events should also have temporal independence in that creating new files and directories should not harm old files and directories”.

The idea of using the filesystem’s recovery tools (fsck) to infer its guarantees seems useful to avoid the requirement of a formally specified model. However, it is our understanding that the resulting model is rather a larger regression test than a truly derived exhaustive model. Such regression tests would be useful as a starting point for a formal model, e.g., as initial litmus tests for use with FERRITE. We do not believe that we can apply the presented approach to ZFS with and without ZIL-PMEM with reasonable effort.

2.1.5 PMEM-specific Crash-Consistency Checkers

A growing body of work introduces tools that check whether code that manipulates persistent memory actually issues the architecturally required instructions for persistence. However, most systems only target userspace code and are thus inapplicable to ZIL-PMEM which is implemented in the ZFS kernel module. Whereas ZIL-PMEM can be compiled for user-space as part of the *libzpool* library, the large amount of conditional compilation involved in the libzpool build process diminishes the significance of user space tests. We surveyed the following userspace-only tools:

- **pmemcheck** [introductionToPmemcheckPart], a tool in Intel’s Persistent Memory Development Kit (PMDK). It integrates with Valgrind ([1]) and requires annotation of all PMEM accesses. Notably, the PMDK libraries include these annotations.
- pmeminsp _____
- Agamotto _____

phrasing is odd?

todo

todo

The remaining tools in this subsection support kernel code and could be applicable to ZIL-PMEM.

Yat: A Validation Framework For Persistent Memory Software [Lan+14]

“Yat is a hypervisor-based framework that supports testing of applications that use Persistent Memory [...] By simulating the characteristics of PM, and integrating an application-specific checker in the framework, Yat enables validation, correctness testing, and debugging of PM software in the presence of power failures and crashes.” The authors used Yat to validate PMFS (see Section 2.1.1). Yat’s hypervisor-based approach makes it the ideal tool for evaluating ZIL-PMEM. To our great dissatisfaction, Yat has never been published and remains an Intel-internal project.

PMTest: A Fast And Flexible Testing Framework For Persistent Memory Programs [Liu+19]

PMTest is a validation tool that claims to be significantly faster than Intel’s *pmemcheck*. The implementation is based on traces of PMEM operations which are generated by (manually or automatically) instrumented application code. PMTest ships with two built-in checkers that assert correct instruction ordering (e.g.: missing store barriers) and durability (e.g.: missing cache flushes). Higher-level checks must be implemented by the programmer. PMTest works within kernel modules but the implementation is limited to a single thread. The processing of the operation trace happens in userspace. PMTest could be used to check ZIL-PMEM kernel code: for a single ZPL filesystem and a single thread that performs synchronous I/O, the pool’s PMEM is only written from that thread.

we introduced
it above, con-
text still there?

ref design sec-
tions

2.1.6 Fault Injection

Error handling code are notoriously difficult to test but often critical for correctness. Fault injection is a common technique to simulate failures and thereby exercise these code paths.

Model-based Failure Analysis Of Journaling File Systems [PAA05] The authors present an analysis of the failure modes caused by incorrect handling of disk write failures in the journaling code in ext4, ReiserFS and IBM JFS. These filesystems each implement one or more of the following journaling modes: data journaling, ordered journaling, and writeback journaling. Any block write performed in one of these modes falls in one of the following categories: “J represent journal writes, D represent data writes, C represent journal commit writes, S represent journal super block writes, K represent checkpoint data writes [...]”. For each of the three journaling modes, the authors present a state machine that describes all permitted sequences of block writes. The state machine includes

transitions to an error state if a block write fails. The system works as follows. A kernel module tracks the filesystems' state as modelled by the state machine. It intercepts block device writes from the filesystem code and injects write failures. If the filesystem subsequently performs another block write that is not permitted by the state machine, an implementation error has been found. The authors distinguish several classes of failures with varying degrees of data loss. The methodology is very filesystem specific which already shows in the adjustments required for IBM JFS. The ZIL's structure and its interaction with txg sync is substantially different from the journaling modes presented in the paper. The system is thus not applicable to ZIL-PMEM.

ref sections

ndctl-inject-error [] is a subcommand of the `ndctl` administrative tool that "can be used to ask the platform to simulate media errors in the NVDIMM address space to aid debugging and development of features related to error handling." The kernel driver forwards injection requests to the NVDIMM firmware via ACPI. [] The errors then surface as *machine check exceptions* (MCE), which is the same mechanism used to indicate detected but uncorrectable ECC errors for DRAM and PMEM. The functionality is useful for testing correct error in the filesystem when *reading* PMEM. ZIL-PMEM does not use `memcpy_mcsafe()` when accessing persistent memory and therefore does not handle MCEs on read. However, if ZIL-PMEM used `memcpy_mcsafe()`, the error handling would be the same as for corrupted or missing log entries, for which we have good unit test coverage.

can we fix this until submission?

ZFS Fault Injection [] The ZFS tool *zinject* allows for targeted injection of artificial IO failures. Errors can be scoped to an entire device or specific logical data objects in ZFS. The ZFS integration test suite makes use of the command for high-level features such as automatic hot spares. ZIL-PMEM does not use ZIO to access PMEM and thus cannot hook into the *zinject* infrastructure. Most of the semantics are tied to ZFS's `zbookmark` and `blkptr` structures which we do not use in ZIL-PMEM. We expect that proper adaptation of *zinject* to ZIL-PMEM would be difficult. A ZIL-PMEM specific tool that allows fault injection or even manipulation of log entries in PMEM would be useful for debugging and integration testing. However, the existing ZIL-PMEM unit tests already test many scenarios for data corruption in a more expressive and efficient manner.

ref evaluation

2.2 Persistent Memory in Linux

2.3 OpenZFS Primer

2.3.1 The ZIL API

reframe this section as dram representation

The *zil.c* code module implements all ZIL functionality. The ZIL is a per-dataset log and thus all ZIL state is kept in per-dataset structures. On disk, this state is kept in the `zil_header_t` structure, which is part of the `objset_phys_t` structure, which in turn is owned by the `dsl_dataset_phys_t` structure. `zil_header_t` contains the pointer to the first LWB in the dataset's ZIL chain. It also contains fields that track claiming and replay progress. In DRAM, the per-dataset state is kept in the `zilog_t` structure. All ZIL APIs are effectively scoped to the `zilog_t`. On the write path, the following ZIL APIs are relevant:

zil_itx_create Allocate an ITX with a given record size.

zil_itx_assign Associate the given ITX with the changes made in the given DMU transaction. The call moves ownership of the ITX to `zilog_t`.

zil_commit Persist previously assigned ITXs to stable storage. The caller can specify the object ID of a `znode` to indicate that it is sufficient to persist the ITXs that affect this `znode`. This functionality enables efficient `fsync()` if some files are written asynchronously and some synchronously.

The activity diagrams in Figure 2.1 visualize how these APIs are used in `write()`, `fsync()`, and `sync()` system calls.

`zilog_t` tracks assigned ITXs in data structures called `itxg`. There exists one `itxg` per unsynced pool transaction group (`txg`). When an ITX is assigned to the ZIL in a given `txg`, it is added to that `txg`'s `itxg`. After the `txg sync` thread has finished syncing a `txg`, it frees the corresponding `itxg` and all the ITXs in it because the changes they describe are now persisted in the main pool and thus obsolete.

Each `itxg` is split into a *sync list* and the *async tree*. The *sync list* is a simple list of ITXs whereas the *async tree* is a search-tree that maps from object ID to a list of ITXs. By default, ITXs are added to the *sync list* when assigned to the ZIL. The *async tree* is only used for ITXs that are scoped to a particular file. For example, an ITX that logs the creation of a file is added to the *sync list* whereas a write within that file is added to the *async tree*.

`zil_commit` uses the `itxgs` to build a linear *commit list* of ITXs that it can subsequently persist to stable storage. The following pseudo-code illustrates the

construction of the commit list by `zil_commit`. Figure 2.2 provides an example for a single transaction group.

```

zil_commit(object_id)
  commit_list := []
  for each unsynced transaction group 'txg':
    itxg <- the itxg for txg
    if object_id == 0:
      append all itxs in itxg.async to itxg.sync
    else:
      append only the itxs in
        itxg.async[object_id] to itxg.sync
      append itxg.sync to commit_list
  err := persist commit_list
  if err:
    wait until the open txg has synced
  return

```

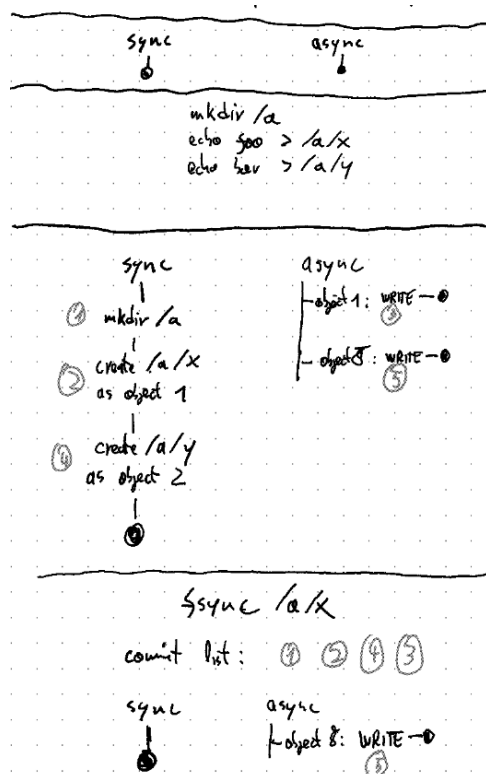


Figure 2.2: An example of how an itxg is filled with ITXs and how `zil_commit` drains it into a commit list. Note that this example only covers the case where all ITXs were assigned for the same txg.

zil_sync

Maybe we can avoid this? It's only relevant in ZIL-PMEM.

2.3.2 Replay

2.3.3 Log Record Types

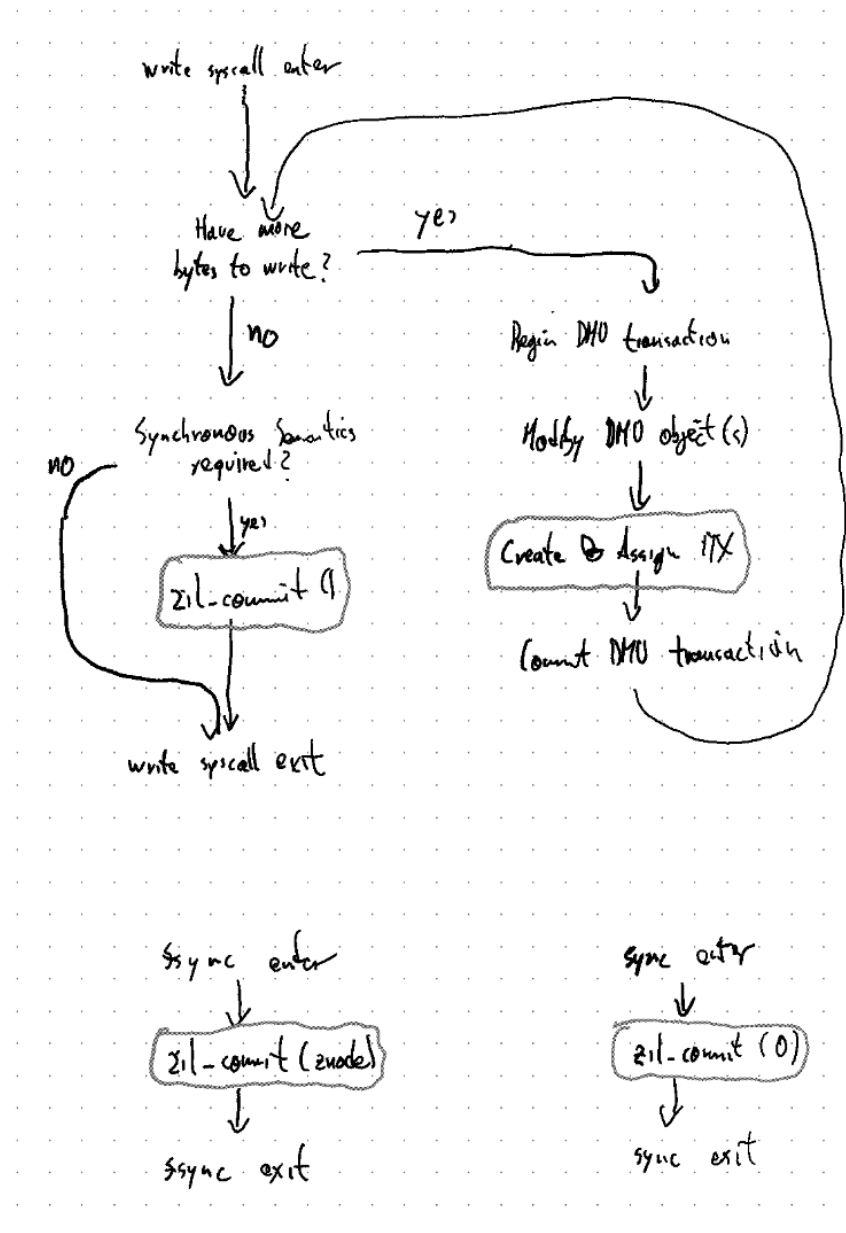


Figure 2.1: ZIL API usage on the write path in `write()`, `fsync()`, and `sync()` system calls.

Chapter 3

ZIL-LWB on PMEM

The motivation for this thesis is the significant overhead of the current ZIL implementation (ZIL-LWB) in 4k random synchronous write workloads compared to the raw PMEM hardware. In this chapter, we describe this benchmark in detail (Section 3.1) and present the resulting data in Section 3.2. We proceed with an analysis of the distribution of the wall clock time spent during this benchmark among the involved ZFS components. Our findings (Section 3.3) show that approximately 80% of overall latency are spent on ZIL-LWB-specific code. We conclude that ZIL-LWB’s data structures and persistence mechanisms are unfit to take advantage of PMEM-level performance, motivating the development of ZIL-PMEM.

3.1 Benchmark Setup

Our benchmarking system has the following hardware configuration:

CPU 2 x Intel(R) Xeon(R) Silver 4215 CPU 2.50GHz
Mainboard Supermicro X11DPi-N(T)/X11DPi-NT, BIOS 3.1a 10/16/2019
DRAM 16 x Micron 8GiB DDR4 2933MT/s (18ASF2G72PDZ-2G9E1), evenly distributed across sockets.
NVMe 3 x Micron PRO 960GB NVMe, 512 byte namespace format (MTFD-HBA960TDF)
PMEM 4 x Intel Optane DC Persistent Memory, 128 GB, (NMA1XXD128GPS), two per socket.

We use the following software stack:

Kernel Linux 5.9, Debian buster (5.9.0-0.bpo.5-amd64)

Userland Debian GNU/Linux (buster)

fio fio-3.23-28-g7064

bpftrace bpftrace v0.12.0

bcc v0.16.0-11-ga74413b0

OpenZFS Our tree of OpenZFS with support for *ZIL kinds* (Section 6.1), configured to use the ZIL-LWB ZIL kind. Note that we ensured manually that the *ZIL kinds* patch do not impact ZIL-LWB performance.

We use the following system configuration:

- We leave SMT enabled, resulting in 16 hardware threads per socket.
- We disable the entire second CPU in software using the `isolcpus=8-15,24-31` kernel command line parameter.
- We configure all Optane DIMMs in *AppDirectNotInterleaved* mode.
- We create a 40 GiB-sized *fsdax* namespace on the region of the first DIMM on the first socket.
- We create a 40 GiB-sized *devdax* namespace on the region of the first DIMM on the first socket.
- We partition the 3 NVMe drives into 10 equals-sized partitions each.
- We create a zpool called `dut` with the 30 partitions as top-level vdevs, and the `/dev/pmem` device as a SLOG. Note that we used the multitude of NVMe partitions because it improved overall device utilization in our setup.
- We create 8 datasets in the zpool, named `dut/ds$i` and mounted at `/dut/ds$i`.
- For all datasets, we configure `recordsize=4k` to match the fio workload and set `compression=off` to avoid CPU overhead in the ZIO pipeline during *txg sync*.

check defined

We use *fio* to generate a workload of random 4KiB writes with synchronous semantics (`blocksize=4k, rw=randwrite`). Each of the one to eight `numjobs` threads performs random synchronous write system calls to a separate file per thread (`ioengine=sync, sync=1, direct=0, fsync=0`). Each file has a size of `size=100MiB` which means that the written data volume grows with `numjobs`, but the amount of dirty data (max.800 MiB) remains well below the NVMe drive's bandwidth limits. To avoid scalability-bottlenecks in the *itxg layer*, we place each thread's file onto a separate dataset `/dut/ds$jobnum/fio_file` using the `filename_format` option. For each value of `numjobs`, we measure for one minute (`time_based=1, runtime=60`), with a ramp-up time of 2 seconds (`ramp_time=2`), and set `end_fsync=1`. Before we run the benchmark, we prepare the benchmark files with a separate *fio* invocation with the additional flag `--create_only=1`.

check this is defined

We apply the *fio* workload to the following storage stacks:

zil-lwb The zpool with ZIL-LWB ZIL kind as described above.

- async** The same configuration as above, but with `sync=disabled` set on all datasets. This turns `zil_commit` into a no-op, reducing the time spent in the write system call to DMU object modification and ITX allocation, all in DRAM.
- fsdax** The fio configuration above, applied directly to the `/dev/pmem` block device (`direct=1` instead of `direct=0`, and `filename=/dev/pmem0` instead of `filename_format=...`). This configuration demonstrates the performance of the Linux block device emulation around the PMEM hardware. It includes the syscall overhead.
- devdax** The fio configuration above, with `ioengine=dev-dax`, applied directly to the `/dev/dax` namespace (`filename=/dev/dax0.1` instead of `filename_format=...`). In this configuration, fio mmmaps the PMEM namespace directly and thus bypasses the kernel completely.

3.2 Results

The following graphs show the achieved IOPS and per-thread latency by numjobs.

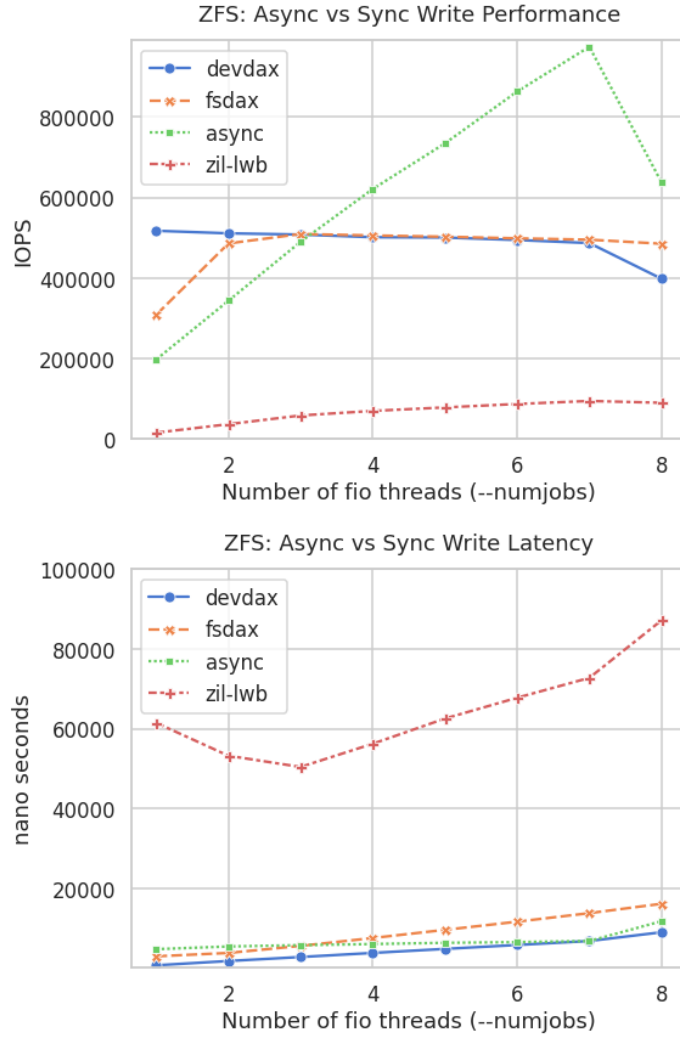


Figure 3.1: IOPS and latency ZIL-LWB compared to ZFS in async mode as well as the raw fsdax and devdax device.

ZIL-LWB only achieves approximately 10k IOPS at one thread and peaks at approximately 100k IOPS at seven threads. In contrast, ZFS in *async* mode starts with 200k IOPS with one thread and achieves over 900k IOPS at seven threads. This exceeds the performance of the PMEM hardware which mostly stays at 500k IOPS and demonstrates that the ZIL is clearly the bottleneck in the *zil-lwb* configuration. A look at the per-IOP average latency emphasizes the vast overhead that ZIL-LWB adds compared to what is possible with the raw PMEM hardware. Whereas ZFS in *async* mode only requires 5 us per IOP with `numjobs=1`, and raw writes to fsdax require approximately 3 us, ZIL-LWB with the same PMEM hard-

ware as SLOG takes more than 60 us per write. The minimum latency achieved by ZIL-LWB is at `numjobs=4` at approximately 50 us before it starts increasing to up to approximately 85 us at `numjobs=8`. Note that we do not use the results from the *dev-dax* as a baseline for PMEM hardware because the IOPS and latencies reported by `fio` do not match. For example, `fio` reports less than *dev-dax* 1 us of latency for *dev-dax* but only reports 550k IOPS at `numjobs=1` whereas $\frac{1}{1 \text{ us/IOPS}} = 10^6$ IOPS would be anticipated at this latency. Figure 3.2 shows the same latency data as Figure 3.1 above, albeit zoomed to a scale that allows us to distinguish *dev-dax*, *fsdax*, and *async* latencies.

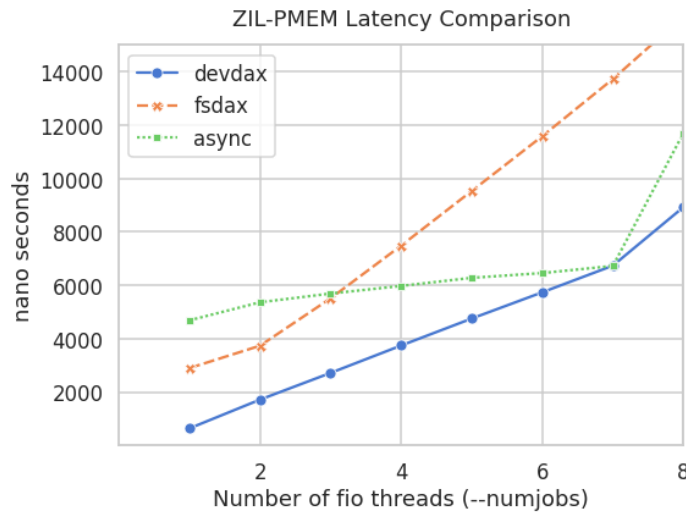


Figure 3.2: Zoomed section of Figure 3.1 that allows us to distinguish *async* mode, *devdax*, and *fsdax*.

3.3 Analysis

By comparing the numbers for *async* and *zil-lwb*, it is safe to assume that regardless of the value for `numjobs`, ZIL-LWB adds at least 40 us of latency. We want to determine where this time is spent.

3.3.1 OpenZFS Background: How ZIL-LWB Works

We have introduced the DRAM representation of the ZIL in Section 2.3.1: Syscalls modify DMU objects in so-called DMU transactions and log the logical change made in a transaction as an ITX to the ZIL. The ZIL stores the ITXs in a data structure called *itxg* until either the transaction group of the ITX's DMU transaction

is synced to disk or synchronous semantics are requested through `zil_commit`. In the latter case, the ZIL constructs the *commit list* which is the sequence of ITXs that needs to be concatenated to some form of per-dataset persistent log structure before `zil_commit` is allowed to return.

ZIL-LWB uses a chain of so-called *log-write blocks* (LWBs) to represent this log structure. The LWB chain is a single-linked list rooted in the ZIL header of each dataset. Each LWB contains a contiguous sequence of variable-length *log records*. A log record is the persistent representation of an ITX. To reconstruct the concatenated commit lists, ZIL-LWB traverses the LWB chain and concatenates the log record sequences within the LWBs. Figure 3.3 visualizes this structure.

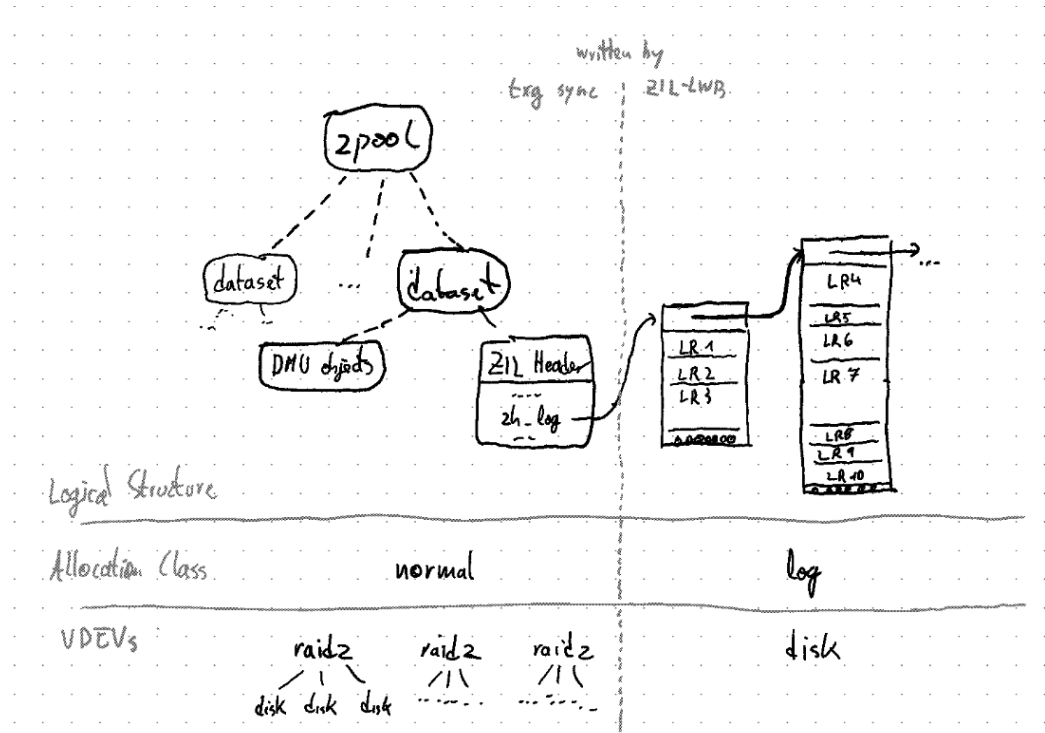


Figure 3.3: On-disk structure of ZIL-LWB as described in the previous paragraph.

The packing of N log records on the commit list into M larger LWBs has been advantageous in the past:

Latency Amortization Under the assumption that disk latency (lat_{disk}) dominates overall synchronous I/O latency, writing log records one-by-one would result in a total latency of $N * lat_{disk}$. In contrast, grouping records into $M \ll N$ LWBs reduces the latency to $M * lat_{disk}$

LWB Timeout / Group Commit ZIL-LWB uses a *timeout* mechanism to extend the amortizing effect of LWB packing. When a thread *A* packs log records into LWBs, it issues the IO operations for all but the last LWB as soon as the LWB is filled. The last LWB remains in *open* state for a short time window. If another thread *B* `zil_commit`s to the same dataset within that time window, it picks up the *open* LWB and continues to fill it with log records from its own *commit list* and issues the LWB IO operation as soon as the LWB is full. The IO completion callback then wakes up *A* and which can now safely return from `zil_commit` because *A*'s commit list is now fully persisted. If the LWB is not committed within the configured time window, *A* issues the IO operation itself. Without the timeout mechanism, thread *B* would need to wait for the last LWB of *A* to be persisted, *and* for its own LWBs to be persisted because *B*'s first LWB is pointed to by *A*'s last LWB. By sharing the last LWB of *A* and *B*, up to $1 * lat_{disk}$ of waiting time can be avoided. Note that whereas ZFS refers to this mechanism as *LWB timeout*, the technique is also well-established in disk-oriented databases under the term *group commit* or *commit group*.

review this,
just thought of
it

Space Efficiency Log records are stored as a contiguous sequence within the LWB. This avoids fragmentation if the commit list consists of small log records because many such log records can be packed into the smallest LWB (4 KiB). However, the on-disk format does not allow for splitting of individual log records. If the log record cannot be split in software (`WR_NEED_COPY`), the open LWB is issued with wasted space and a new LWB is allocated for the log record.

LWBs are allocated from the metaslab allocator which identifies the location of allocated data through a *ZFS block pointer*. Block pointers are fat pointers (128 byte) that, among other metadata, contain up to three *data virtual addresses* (DVAs) that identify where one to three copies of the data are stored. A DVA consists of the top-level vdev ID and an offset in the vdev's logical block address space — other systems would probably refer to this as logical extent(s). Block pointers also store the checksum of the data that they point to. Computing this checksum is possible because ZFS's main on-disk structure is a tree that is written bottom-up by the *txg sync* thread for every transaction group. The end result is a merkle tree that is rooted in a block pointer stored in ZFS's *uberblock* which is roughly equivalent to superblocks in other filesystems.

The LWB chain is a special case because the ZIL must be able to add new LWBs between transaction groups. The solution is to pre-allocate LWBs and to store the checksum of their content in the LWB itself. The procedure to write out an LWB thus is as follows:

1. Wait until the LWB is filled or the timeout triggers the LWB to be written out.
2. Predict the best size for the next LWB based on a simple heuristic.
3. Allocate the next LWB, preferably from the *log* allocation class.
4. Store the resulting block pointer in the current LWB so that it points to the next LWB.
5. Compute the checksum of the current LWB.
6. Repurpose the current LWB's block pointer field to store the computed checksum.
7. Write out the current LWB.

Figure 3.4 visualizes these steps for the initial case where no log exists, and for an existing log. The LWB chain's end is implicitly marked by an invalid checksum. The checksum also ensures crash-consistency since the list is append-only and a partially written block is equivalent to the end of the log during recovery. Note that it is unlikely for an unwritten but allocated LWB to be mistaken for a valid LWB block due to additional metadata stored in the LWB. For OpenZFS native encryption, the LWBs are not checksummed but encrypted and authenticated using an AEAD algorithm such as aes-256-gcm.

I think cipher mode is not the right term

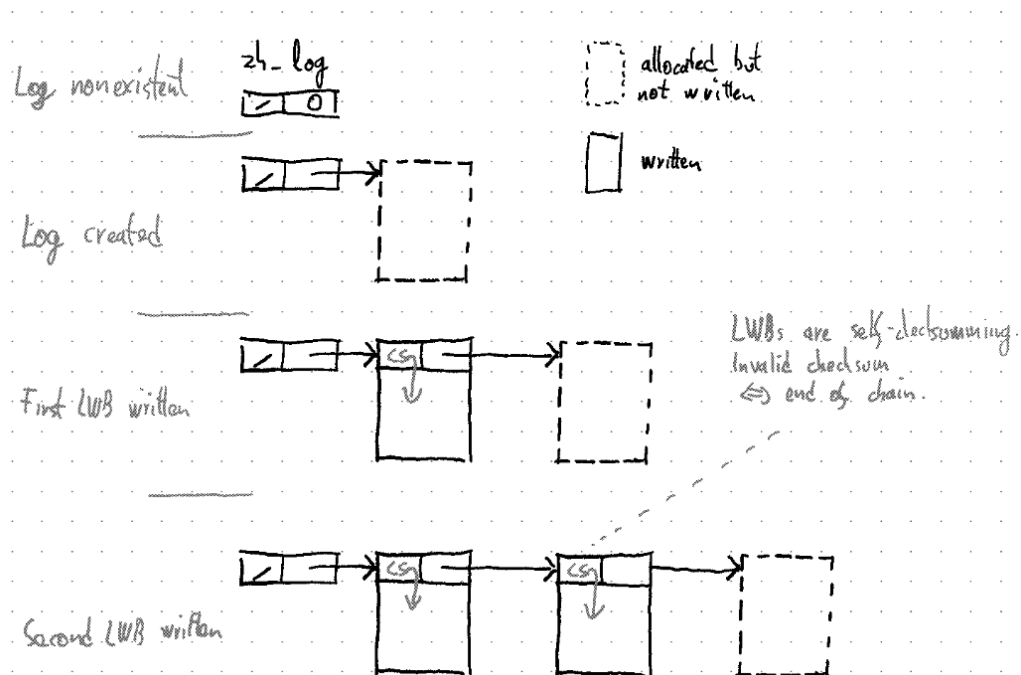


Figure 3.4: The procedure for appending an LWB.

ZIL-LWB does not write LWBs directly using the kernel's block device API but uses the ZFS *ZIO pipeline* instead. ZIO is ZFS's unified abstraction for performing deferred allocation, transparent data and metadata checksumming, compression, deduplication and encryption as well as redundancy mechanisms such as raidz. The architecture is pipeline-oriented and heavily parallelized using ZFS's *taskqs* to maximize throughput for its main consumer, *txg sync*. By using the same infrastructure for LWB ZIO operations, a very high degree of code reuse is possible, although there are some special cases, in particular for self-checksumming or encryption. Also, integrating LWB I/O into the ZIO pipeline could be beneficial for the case where the pool does not have a SLOG configured because latency-sensitive ZIO operations can be prioritized over throughput-oriented *txg sync* operations. However, as we will elaborate on in the next subsections, ZIO adds significant latency overheads, resulting in sub-par performance of ZFS on low-latency storage devices, and PMEM in particular.

Analysis

Given the background knowledge on how ZIL-LWB works, we are now ready to present our latency analysis. We use dynamic instrumentation (eBPF via bpftrace) to sum up the wall clock time spent in ZFS functions that are executed by the fio threads during a synchronous write. We then use the model in Figure 3.5 to compute the time spent in the asynchronous part of the write operation, the ITX layer, and ZIL-LWB specific code. We visualize the results as stacked bar charts in Figure 3.6.

cite

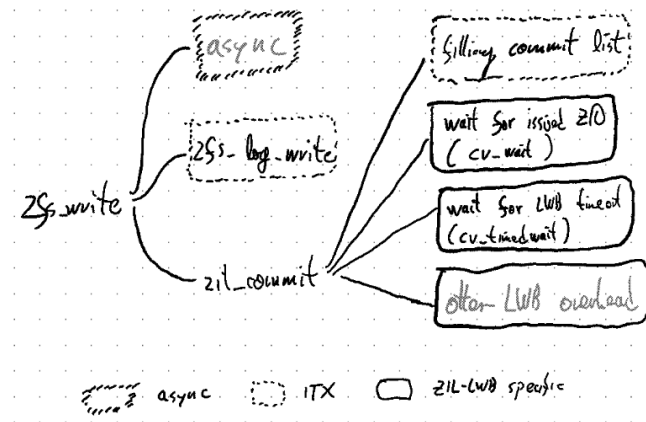


Figure 3.5: Our model of the time spent in ZFS by a fio thread that performs a write system call. The different columns represent the different levels of the dynamic call graph. The nodes in each column describe all activity that happens at this level of the call graph. Activities with black text color are instrumented using eBPF. Gray text color indicates that the value was computed in post-processing by subtracting the sum of the columns instrumented time from the instrumented time of the parent in the column to the left. For example, to compute *async*, we compute $zfs_write - (zil_commit + zfs_log_write)$. The border style of the node visualizes the subsystem to which we attribute the time spent in the activity represented by the node.

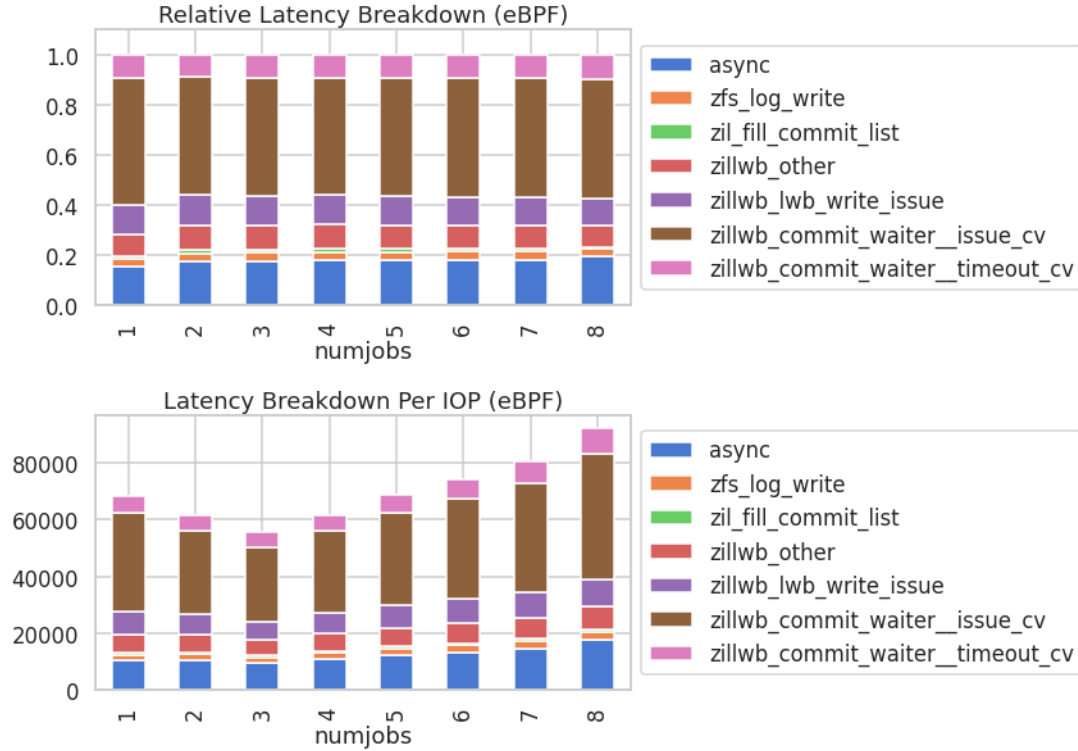


Figure 3.6: The relative and absolute breakdown of latency by activity. The absolute breakdown is normalized by the number of write operations. We compared the latencies in the breakdown with the latencies observed by fio to ensure that our accounting is correct under the given model.

Our observations are as follows:

- The instrumentation overhead is 5–10 us per IOP. (We compared the latencies that fio reports for the instrumented and uninstrumented run.)
- The relative distribution of latency remains mostly unchanged for the different values of numjobs.
- The *async* part of the write path and the ITX layer together only amount to approximately 20% of overall latency. The remaining 80% are spent on LWB-specific activities.
- The LWB timeout mechanism amounts to 4% of overall latency.
- At least 20% of overall latency are spent on filling LWBs, issuing their corresponding ZIOs, and other LWB-related activities.
- 45–50% of overall latency is spent waiting for the ZIO pipeline to persist the LWBs. In absolute numbers, the value ranges from 25–40 us. We have separately confirmed that ZIL-LWB uses 12 KiB LWBs which is the smallest

possible LWB size allowed by the implementation. In our benchmarking setup (per-thread datasets), each LWB only holds a single 4 KiB write log record with 192 B of additional metadata. This amounts to write amplification of 3x. However, even with this increased data volume per IOP, the results from the previous experiment suggest that with a conservative estimate of 3 us per 4k write to the raw PMEM hardware, we should expect less than 9 us of PMEM write time per LWB. Thus, a major fraction of the 25–40 us that are spent waiting for ZIO is actually pure software overhead.

Our observations lead us to the following conclusions:

First, we have reason to believe that the high overhead of ZIO is unlikely to be reduced to a degree that allows for exploitation of PMEM-level latency, regardless of whether the hardware is actual PMEM or simply very fast NVMe drives. The reason is that there is an inherent conflict of goals for ZIO: whereas the ZIL is strictly latency-oriented, all other consumers (*txg sync*, *scrubbing*, *zfs send/recv*) are throughput-oriented. The pipeline-oriented, parallelized architecture of ZIO is important for throughput and provides great flexibility, but increases latency through context switches. This problem is well known in the ZFS community, see [\[Ope20\]](#).

more quotes,
find gh issues,
etc

Second, the persistent representation of the ZIL as a chain of LWBs poses a severe and unnecessary overhead on PMEM. For one, the latency amortization provided by LWBs and the timeout mechanism is unnecessary at PMEM-level latencies where it is cheaper to persist the log records on an individual basis than batching them in LWBs and coordinating with other threads on the matter. And for another, since PMEM is byte-addressable, the persistent representation is no longer constrained by disk block sizes, opening up the possibility of better space efficiency.

Chapter 4

Design Overview

Given the insights described in the previous chapter, we propose an alternative ZIL implementation called **ZIL-PMEM** that exclusively targets PMEM SLOG devices. In this chapter we define the project goals for ZIL-PMEM and provide a high-level overview of its design. The subsequent two chapters then introduce the core data structure and our approach to integrating it into ZFS.

4.1 Project Goals & Scope

4.1.1 Requirements

Coexistence ZIL-PMEM must coexist with ZIL-LWB in code and at runtime due to limited availability of PMEM hardware and the limitations of the ZIL-PMEM design.

Same Guarantees ZIL-PMEM must maintain the same crash consistency guarantees towards user-space as ZIL-LWB for both ZPL and ZVOL.

Simple Administration & Pooled Storage Pooling of storage resources and simple administration are central to ZFS [Bon+03]. ZFS should automatically detect that a SLOG device is PMEM and, if so, use ZIL-PMEM for all of the pool's datasets. No further administrative action should be required to fully benefit from ZIL-PMEM.

Correctness In the absence of PMEM media errors and data corruption, ZIL-PMEM must be able to replay all data that it reported as committed. The result must be the same as if ZIL-LWB would have been used in lieu. Specifically:

- Replay must respect the logical dependencies between log records.
- Logging must be crash-consistent, i.e., the in-PMEM state must always be such that replay is correct.
- Replay must be crash-consistent, i.e., if the system crashes or loses power during replay, it must be possible to resume replay after the crash. Resumed replay must continue to respect logical dependencies of log records.

Data Integrity Data integrity is a core feature of ZFS [Bon+03]. ZIL-PMEM must detect corrupted log records using an error-detecting code. Detected corruption must be handled *correctly* (as outlined in the previous paragraph) and *gracefully* with the following behavior as the baseline: “Assume a sequence of log records $1 \dots N$ where log record 1 does not depend on a log record and each record $i > 1$ depends on its predecessor $i - 1$. Data corruption in record $i \in 1 \dots N$ must not prevent replay of records $1 \dots i - 1$ ”.

Low Latency The latency overhead of ZIL-PMEM compared to raw PMEM device latency for the same data volume should be minimal for single-threaded workloads. Multi-threaded workloads are addressed below.

Multi-Core Scalability Since PMEM is added as a pool-wide resource used by all of the pool’s datasets, ZIL-PMEM should scale well to multiple cores. Barring PMEM throughput limitations, the speedup in throughput (IOPS) achieved by parallelizing synchronous I/O to multiple cores on a ZIL-PMEM system should be as follows:

1 private dataset per thread Always near-linear speedup.

1 shared dataset

ZPL filesystem No speedup.

ZVOL No speedup in standard mode, potentially sub-linear speedup in bypass mode (Section 6.3.5).

Maximum Performance On Intel Optane DC Persistent Memory We develop and evaluate ZIL-PMEM exclusively for/on Intel Optane DC Persistent Memory since it is the only broadly available non-volatile main memory product on the market. Whereas supercapacitor-backed persistent memory modules (NVDIMM-N) should be usable with ZIL-PMEM, our goal is to design a system that makes optimal use of the Optane hardware.

CPU-Efficient Handling Of PMEM Bandwidth Limits If the maximum write bandwidth to any kind of storage device is exceeded, the I/O stack must somehow apply back-pressure to avoid losing in-flight data. With PMEM, the I/O

stack is the CPU microarchitecture and the back-pressure manifests as stalling instructions. However, from the OS thread scheduler's perspective, threads whose instructions stall because they wait for PMEM are indistinguishable from actually busy threads. Yang et al. have shown that a single Optane DIMM's write bandwidth can be exhausted by one CPU core at 2 GB/s and that write bandwidth decreases to 1 GB/s at ten or more CPU cores. Since ZIL-PMEM shares PMEM among all datasets in a zpool, we expect bandwidth exhaustion to be a phenomenon that will happen in practice. ZIL-PMEM should thus provide a mechanism to shift excessive PMEM I/O wait time off the CPU.

Testability ZIL-PMEM must be architected for testability. The core algorithms must be covered by unit tests. Further, ZIL-PMEM should be integrated into the ztest user-space stress test as well as the SLOG tests of the ZFS Test Suite.

4.1.2 Out Of Scope For The Thesis

The following features were omitted to constrain the scope of the thesis. We believe that our design can accommodate them without major changes.

Support For OpenZFS Native Encryption The ZIL-PMEM design presented in this section does not address OpenZFS native encryption. Intel Optane DC Persistent Memory supports transparent hardware encryption per DIMM at zero overhead. In contrast, OpenZFS native encryption is per dataset and software-based. Given these significant differences in data and threat model, ZIL-PMEM cannot rely on Optane hardware encryption. Instead, ZIL-PMEM would need to invoke OpenZFS native encryption and decryption routines when writing or replaying log entries.

cite spec

Protection Against Scribbles Scribbles are bugs in the system that accidentally overwrite PMEM, e.g., due to incorrect address calculation or out-of-bounds access in the kernel. PMEM-specific filesystems such as PMFS and NOVA-Fortis have already introduced mechanisms to protect against scribbles [Dul+14; Xu+17]. We believe that these mechanism can be applied to our design as well.

4.1.3 Limitations

The following features are deliberately not addressed by our design. More experimentation and experience with ZIL-PMEM will be necessary to determine which features are useful in practice, how they can be realized, and how they interact with the existing requirements.

No NUMA Awareness Yang et al. recommend to “avoid mixed or multi-threaded accesses to remote NUMA nodes. [...] For writes, remote Optane's latency is

2.53x (ntstore) and 1.68x higher compared to local" [Yan+20]. We do not account for this behavior in the design and do not evaluate ZIL-PMEM in a NUMA configuration.

future work

No Data Redundancy ZIL-PMEM provides data integrity protections but does not provide a mechanism for data redundancy.

future work

Only Works With SLOGs Our approach to integrate ZIL-PMEM into ZFS is only applicable to PMEM SLOGs and does not work for a zpool that uses PMEM as main pool vdevs. Such pools continue to use ZIL-LWB.

No Software Striping Our design only supports a single PMEM SLOG device. Users may wish to use multiple PMEM DIMMs to increase log write bandwidth. With Intel Optane DC Persistent Memory, multiple PMEM DIMMs can be interleaved in hardware with near-linear speedup [Yan+20]. Whereas software striping would be the natural approach to ZFS, it will be non-trivial to achieve the same speedup as hardware-based interleaving.

No Support For WR_INDIRECT ZIL-LWB writes the data portion of large write log records directly to the main pool devices. The ZIL record then only contains metadata such as `mtime` and a block pointer to the location in the main pool. This technique avoids double-writes which is particularly advantageous if the pool does not have a SLOG, which in turn is a use case that ZIL-PMEM does not address (see above). Further, if a SLOG is available, `WR_INDIRECT` log record write latency is likely to be dominated by the main pool's IO latency if it consists of regular block devices. If the main pool's IO latency were acceptable, a fast NVMe-based ZIL-LWB SLOG or no SLOG at all would likely be sufficient for the setup in question.

Space Efficiency ZIL-PMEM is allowed to trade PMEM space for time and simplicity when presented with the option. Our justification is twofold. First, PMEM capacities are significantly higher than DRAM. For example, the smallest Intel Optane DC Persistent Memory DIMM offered by Intel is 128 GiB. [`optanepricing_missing`] Second, the maximum amount of log entry space required from any ZIL implementation is a function of the maximum amount of dirty data allowed in the zpool. For ZIL-LWB, small SLOG devices of 16 to 32 GiB are sufficient in practice. Thus, there is sufficient headroom for PMEM space usage in ZIL-PMEM.

ref background section? this sentence should remain, though

4.2 Design Overview

ref ix systems truenas M

We introduce the concept of *ZIL kinds* to ZFS. The ZIL kind is a pool-scoped variable that determines the pool's strategy for persisting ZIL entries. A zpool's

ZIL kind is determined by the following rule: if the pool has exactly one SLOG and that SLOG is PMEM, the ZIL kind is ZIL-PMEM. Otherwise, it is ZIL-LWB. As explained in Section 3.3.1, ZIL-LWB uses ZFS's metaslab allocator to allocate log-write blocks (LWBs) from the storage pool with a bias towards SLOG devices. In contrast, ZIL-PMEM disables metaslab allocation for the PMEM SLOG device and uses the PMEM space directly.

We partition the PMEM SLOG's space into fixed-size contiguous segments called *chunks*. We develop a data structure called *PRB* that consumes these chunks and exposes the abstraction of an unordered persistent storage layer for *log entries*. A log entry is the unit of data that can be written to and read from PRB. It consists of the ZIL's log *record* and PRB-specific metadata. PRB scales to many concurrent writers and features a mechanism to avoid excessive on-CPU waiting for PMEM I/O. Log entries stored in PRB are automatically garbage-collected when they become obsolete because their transaction group has been synced.

PRB provides a pool-wide storage substrate for log entries but does not define any structure. This is the role of the *HDL* abstraction which implements a mostly sequential log structure on top of PRB. Each dataset in the zpool has a separate HDL to which it writes log entries. The HDL adds metadata to the entry that attributes it to the HDL and encodes the log's structure. After a system crash, the HDLs scan PRB for entries that need to be replayed and *claim* them to hold them back from garbage collection. For replay, the HDL API provides a callback-based interface that allows its consumer to apply the changes that are encoded in the log entries in sequential, deterministic order. Loss of entries, e.g., due to data corruption, is handled as gracefully as possible under the constraints of the logical dependencies encoded in the log structure.

Pools of the ZIL-PMEM ZIL kind use PRB/HDL to persist ZIL log records. For this purpose, we refactor the existing *zil.c* code module. Before the introduction of ZIL kinds, the struct *zilog_t* implemented all ZIL-related functionality. With ZIL kinds, *zilog_t* acts as an abstract base class that encapsulates shared ZIL functionality. This includes the definitions of the ZIL log record format and the *itxg* data structure that tracks which log entries need to be persisted when synchronous semantics are requested by a syscall. The code that is responsible for persisting log entries resides in ZIL-kind-specific subclasses. The pool's ZIL kind determines which subclass is instantiated at runtime. The subclass for ZIL-LWB is *zilog_lwb_t* which contains the original LWB code that we move from *zil.c* into the new *zil_lwb.c* module. The subclass for ZIL-PMEM is *zilog_pmem_t*. It is a thin wrapper around HDL's methods for writing and replaying log entries. *zilog_pmem_t* is implemented in the *zil_pmem.c* code module. This module also

language

contains the code that sets up PRB/HDL on top of the PMEM SLOG vdev, and the code that synchronizes the lifecycles of HDLs their corresponding datasets.

check

The next two chapters describe our design in detail. Chapter 5 describes our main contribution — the PRB/HDL data structure. The integration of PRB/HDL into ZFS is then presented in Chapter 6.

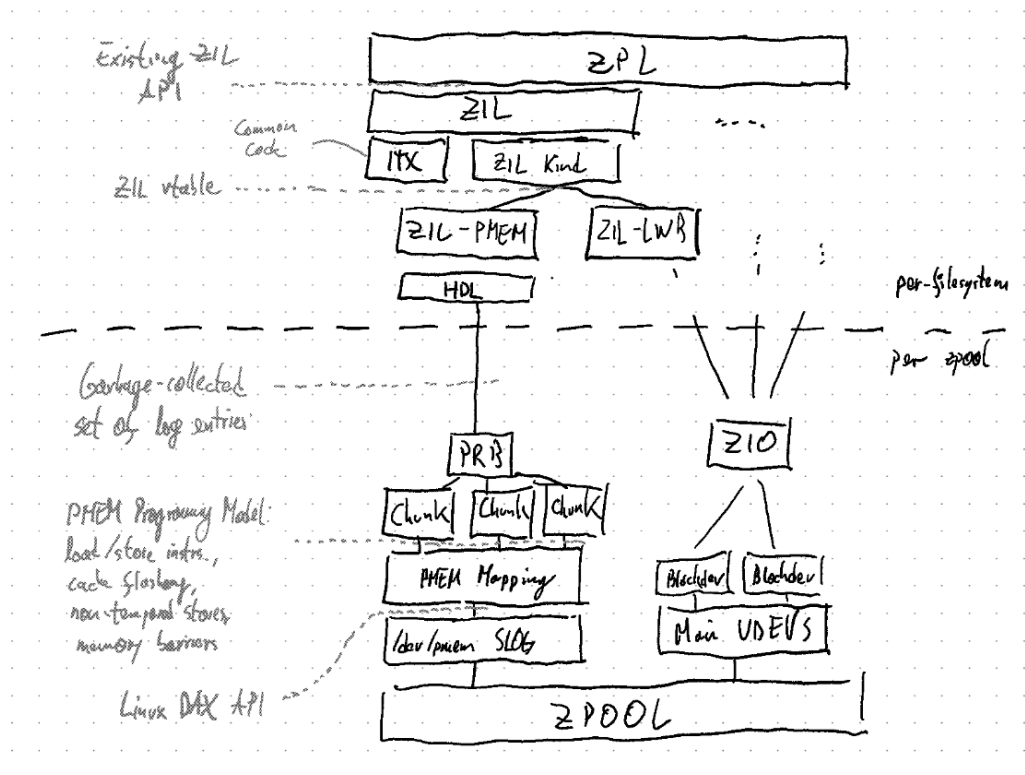


Figure 4.1: Overview of the system architecture as described in this section.

actual figure

Chapter 5

The PRB/HDL Data Structure

In this chapter we describe the PRB/HDL data structure which implements the bulk of ZIL-PMEM's functionality. PRB/HDL abstracts the allocatable space of a PMEM SLOG vdev into virtual logs for each dataset in the pool. The *zil_pmem.c* module, which we present in Chapter 6 uses these virtual logs to implement the ZIL-PMEM ZIL kind. We present the design and implementation of PRB in a top-down manner. In Section 5.1 we recapitulate the role of the ZIL in ZFS to subsequently analyze the requirements that ZIL-PMEM puts on PRB in Section 5.2. Afterwards, in Section 5.3, we give a high-level overview of our design. Section 5.4 presents the virtual log abstraction that is exposed by HDL and Section 5.5 describes the high-level approach for replay. Sections 5.6, 5.7, and 5.8 then progressively refine our understanding of replay and explain how data corruption and crash-consistency is addressed by HDL. Subsequently, we describe PRB which is the storage substrate that the HDLs use for persistence: Sections 5.9 and 5.10 present the data structures that we use for PMEM space management and log entry storage. In Section 5.11 we describe the algorithm that traverses the in-PMEM data structure during log recovery. Section 5.12 then explains how garbage collection removes entries from it. The low-latency and CPU efficient design of the write path is then presented in Section 5.13. Finally, Section 5.14 provides an overview of the PRB/HDL API that is consumed by *zil_pmem.c*.

language, help

help, don't want to repeat 'in-PMEM structure'

5.1 OpenZFS Background

Remember from Section ?? that whenever a file system call changes a dataset D , it does so in a DMU transaction T_i within a transaction group T_{itxg} . After the system call handler has finished the DMU transaction by calling `dmu_tx_commit(T_i)`, the logical change C_i made in T_i is not yet persisted to stable storage. Instead,

the DMU accumulates the changes from many DMU transactions in DRAM as so-called *dirty state*, grouped by the transaction's txg. Eventually, the *txg sync* background thread syncs out a new version of the zpool state that contains the accumulated changes of the txg: it first *quiesces* the currently *open txg* by not admitting new DMU transactions to it and waits for existing DMU transactions to finish. Once all DMU transactions have finished, it is guaranteed that the accumulated dirty state for this txg is not going to change. The txg transitions to the *syncing* state and txg sync starts to write out an updated version of the on-disk state. After this process is complete, the transaction group is the pool's new *last synced* txg. Note that there are three unsynced txgs at any given time, one for each of the three states *open*, *quiescing*, and *syncing*. This improves DMU transaction latency because it decouples new DMU transactions from the *txg sync* thread: new DMU transactions can always operate on the open txg while txg sync only operates on the syncing txg.

The ZIL bridges the gap between the time at which a DMU transaction T_i is finished (`dmu_tx_commit`) and the time at which the changes made by T_i actually reach the main pool through *txg sync*. For that purpose, the system call handler encodes the change C_i that was made in T_i as a *log record*, wraps it in an ITX and assigns it to the ZIL. The ZIL stores assigned ITXs in DRAM in the *itxg* data structure. If the system call has synchronous semantics, the system call handler invokes `zil_commit` before returning to userspace. `zil_commit` uses *itxg* to determine the *commit list* which contains the sequence of ITXs that needs to be concatenated to some form of persistent log structure. In upstream ZFS, the ZIL(-LWB) implements this persistent log structure by stuffing ITXs into LWBs and linking the LWBs into a chain on disk. With the introduction of ZIL kinds, `zil_commit` only builds the commit list and leaves the persistence to the ZIL kind implementation.

After a system crash, the zpool is in the state of the last synced transaction group which we call *precrash_txg*. Our dataset D is lacking exactly those changes C_i whose transactions T_i were made in transaction groups that are younger than *precrash_txg*, i.e., their $T_{itxg} > \text{precrash_txg}$. When D is mounted, these changes C_i need to be applied to D in order to recover data that was reported as committed towards userspace. The ZIL API for this task is called `zil_replay`. Its implementation is specific to each ZIL kind but the API contract is the same: the caller provides a callback that the implementation invokes for every log record that represents a missing change C_i . The callback interprets the log record and performs a DMU transaction that atomically applies the change to the dataset and records replay progress in the ZIL header. The invocation order is defined by *itxg* through the commit list order at write time, but entries for transactions in txgs $T_{itxg} \leq \text{precrash_txg}$ must be skipped because the change encoded in

them is already part of the dataset state. We provide an example for the commit list and the replay sequence for different precrash-txgs in Figure 5.1.

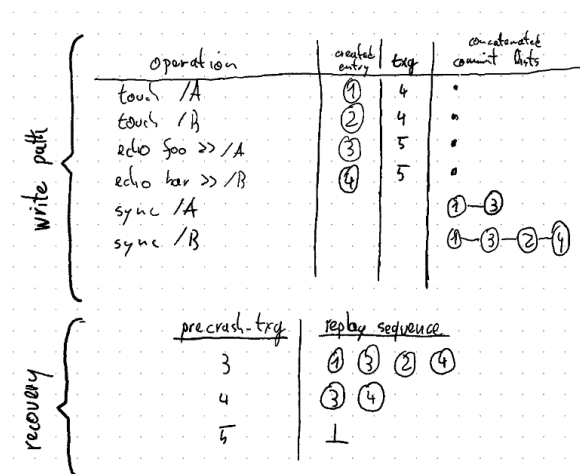


Figure 5.1: Example for a commit list with entries from different files and txgs. We show the replay sequences for different precrash-txgs. The replay callback must only be invoked for the entries that log changes in txgs younger than the precrash-txg.

Note that the DMU transactions during replay only re-enact the logical changes to the dataset but do not restore the exact physical state of the dataset that would have been synced in the absence of a crash. For example, changes that would have been spread across transaction groups 23, 24, and 25 at write time may land in transaction groups 42 and 43 during replay. Further, the system could crash again, during replay. For example, we could lose power so that txg 42 is the last-synced txg whereas 43 was still in *syncing* state. Since the replay callbacks for the individual log records are not idempotent, the ZIL must ensure that, after the crash during replay, the changes that were already replayed in transaction group 42 are not replayed again. For example, the second replay attempt could apply the remaining set of changes in transaction group 83, but it may also be possible that they are spread over several txgs, e.g., 83 and 84. The consequence for any ZIL implementation is that it needs to persist the following data, e.g., in the ZIL header:

Precrash-txg The *precrash-txg* for open logs at the first time during pool import so that it can filter log entries by that criterion, and

Replay progress Information that tracks which log entries have already been replayed so that replay can be resumed correctly after a crash.

5.2 Requirements

The *commit list* is the point at which the common ZIL code hands over responsibility to the ZIL kind implementation. We identify the following requirements for the persistence layer of any ZIL kind:

- During normal operation, it acts as a write-only sink for ITXs whose order is defined by the *itxg* data structure in the form of the *commit lists*.
- The persistent representation of an ITX is the *log record*. A log record is the encoded representation of the change that was made in a single DMU transaction. Log records are always scoped to a single dataset. Their representation is shared among all ZIL kinds. However, the ZIL kind introduces additional structures to encode the commit list order. For example, ZIL-LWB introduces the log-write Block (LWB) as explained in Section 3.3.1. ZIL-PMEM wraps log records in so-called *log entries* that contain additional metadata.
- In the event of a crash, the persistence layer must replay exactly those log records that have been successfully written to the ZIL but whose DMU transaction's *txg* did not sync before the crash. The code that decodes the log records and applies the changes encoded in them is shared among all ZIL kinds. It applies each log record's change in separate DMU transactions and gives the persistence layer an opportunity to update the ZIL header in the same transaction.
- Each dataset has a separate log that is written and replayed independently.

We derive the following abstract view of **what** needs to be stored **per log**:

- The log records themselves.
- The transaction group of the DMU transaction that the log record encodes.
- Structural information that defines replay order and/or logical dependencies between log entries that replay must respect.
- The *precrash-txg* to discern replayable from obsolete log records (see previous section).
- Some representation of *replay progress* to enable resumption of replay if the system crashes during replay.

For ZIL-PMEM, we put the following requirements on the **storage substrate** that stores the log entries:

- On the write path, the overhead added to the raw log entry write time should be minimal.
- It must scale well on a multicore system since many datasets write their log entries in parallel. This scalability requirement includes efficient use of CPU time in case the PMEM write bandwidth is exceeded.
- The storage substrate is responsible for garbage-collecting log entries after they are obsolete, either by txg sync during normal operation or because they have been replayed.
- It must provide a facility to retrieve non-obsolete log entries of a dataset for replay.
- It must detect data corruption using checksums. (Repair and redundancy are out of scope for this thesis, see Section 4.1.1.)

5.3 Approach

We introduce the pool-wide *PRB* object which abstracts the PMEM SLOG vdev's space as a persistent, unordered set of log entries. A log entry encodes a logical change to a particular dataset that was applied in a single DMU transaction. PRB provides facilities for adding log entries to the set and for iterating over its contents. It automatically garbage-collects entries after they become obsolete because their transaction group has synced.

Log entries are created by the thread that performs the DMU transaction. However, the thread does not write them directly to PRB but to the *HDL* object of the modified dataset. A HDL is a virtual log built on top of PRB that organizes individual entries for a single dataset in a mostly sequential structure. After a system crash, the HDL provides a replay facility that recovers the replayable entries, orders them according to their logical dependencies, and handles missing entries.

At any time, there exists one HDL for each head dataset in the pool. They are set up early during pool import and torn down late during export. If a dataset is created or destroyed, the corresponding HDL is set up or torn down as well. HDL is stateful. Its internal states represent the different phases that a dataset goes through with regards to the ZIL.

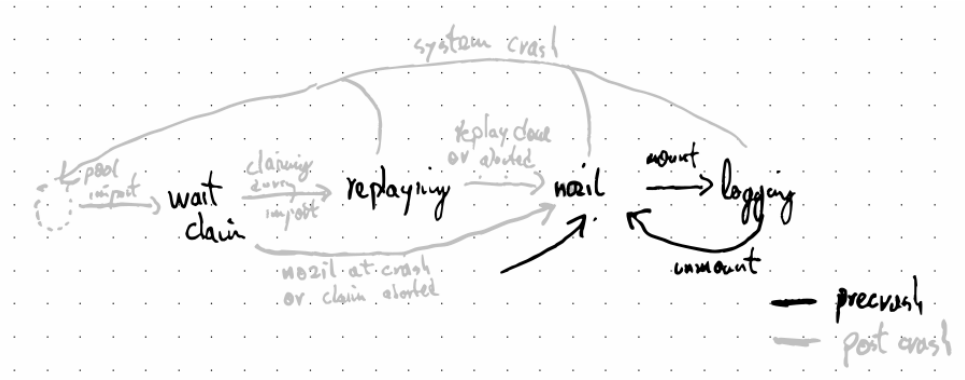


Figure 5.2: TODO

When the HDL is created, it does not have a log and is in state *nozil*. When the dataset is mounted, the HDL allocates a log GUID which uniquely identifies the log's entries in the PRB and persists it to the ZIL header. The HDL is now in state *logging* and threads can write entries to it. If the dataset is unmounted, the log GUID is discarded and the log transitions back to state *nozil*. Otherwise, the HDL remains in state *logging* until the system crashes.

When a zpool is imported, the *claiming* phase examines the ZIL header of each dataset to recover the HDL's runtime state from before the crash. If the HDL state was *logging*, the claiming procedure saves the zpool's *last synced* txg as the *precrash-txg* and transitions the HDL and ZIL header to state *replaying* in initial replay position. If the HDL was already in state *replaying*, the precrash-txg and replay position are recovered from the ZIL header. At this point, all HDLs are either in state *nozil* or *replaying*. HDLs in state *replaying* scan the PRB for log entries that need to be held back for replay. Entries that are held back by at least one HDL are exempt from PRB's garbage collection. After this step, the claiming phase is done and the txg sync thread starts. HDLs are not replayed until their dataset is mounted which happens on a per-dataset basis at a user-controlled point in time. After replay is complete, the HDL discards the log GUID and transitions to state *nozil*. At this point, the mount procedure behaves as if log replay had not happened and starts a new log with a new log GUID, thereby closing the circle. Note that the HDL (and PRB) are able to tolerate a system crash in any of the HDL or ZIL header states. We address this issue in Section 5.8.

5.4 HDL: Log Structure

The structure of the virtual log that each HDL represents is defined by metadata stored with each entry that is written to PRB.

We **attribute** entries to a given HDL's log through the *log GUID*:

Log GUID A 128 bit random identifier stored in the HDL's ZIL header and repeated in every entry written through that HDL.

The following pieces of metadata define the structure of the log.

Transaction Group (txg) The transaction group in which the change encoded in the entry was or would have been synced out by *txg sync*.

Generation Number (gen) The log is a sequence of generations, each of which contains many log entries. The generations encode logical dependencies between entries. Entries within the same generation do not depend on each other. Entries from newer generations unconditionally depend on all entries in all previous generations. We represent generations as unsigned 64-bit non-zero integers.

Generation-Scoped ID (gsid) Within a generation, we identify entries by another by the *gsid*, another unique unsigned 64-bit non-zero integer. As the name suggests, the *gsid* only needs to be unique within a generation.

Note that the tuple $(gen, gsid)$ uniquely identifies an entry within a log.

We visualize the structure of the log in a grid. The rows represent the transaction group (*txg*) and the columns represent the generation (*gen*). For readability, we represent entries not by $(gen, gsid)$ but by a single unique letter. The projection of entries onto the horizontal axis shows the dependency relationship encoded by the generations. The projection of entries onto the vertical axis shows the sets in which entries are garbage-collected.

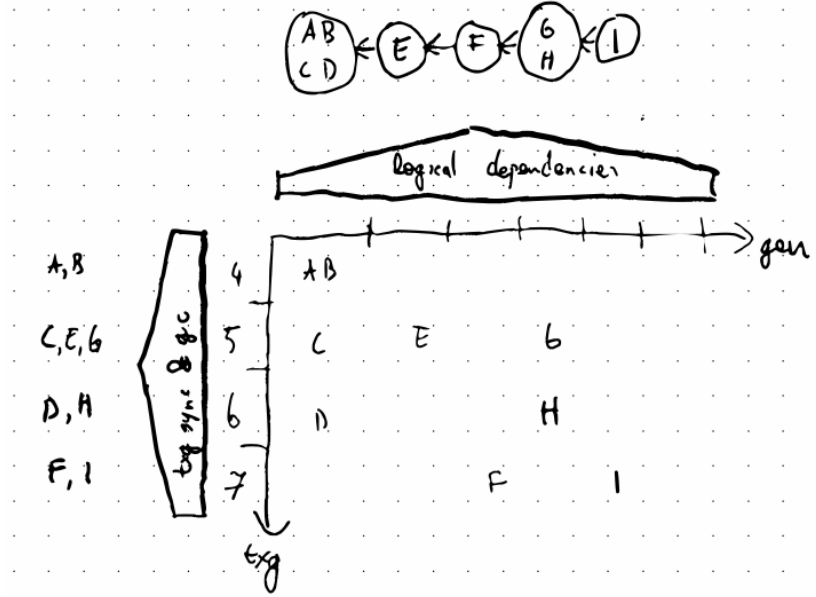


Figure 5.3: Structure of a single HDL's log as described in the previous paragraph.

5.5 HDL Replay: Basic Approach

Replay must apply the changes that were successfully logged to the HDL but whose DMU transaction did not sync before the system crashed. To accomplish this task, it scans the PRB for entries that belong to the HDL and determines a *replay sequence* S .

$$\begin{aligned} &\text{Entry } E_i \in S \\ \Leftrightarrow &E_i.\text{log_guid} = \text{HDL.log_guid} \wedge E_i.\text{txg} > \text{precrash_txg} \end{aligned}$$

S is sorted in *replay order*, which is the order the changes encoded in the entries need to be applied. Given two entries $E_a, E_b \in S$ for the same HDL, the replay order is a total order defined by

$$\begin{aligned} &E_a <_{\text{replay}} E_b \\ \Leftrightarrow &(E_a.\text{gen}, E_a.\text{gsid}) <_{\text{lexicographical}} (E_b.\text{gen}, E_b.\text{gsid}) \\ \Leftrightarrow &E_a.\text{gen} < E_b.\text{gen} \vee (E_a.\text{gen} = E_b.\text{gen} \wedge E_a.\text{gsid} < E_b.\text{gsid}) \end{aligned}$$

For our example in Figure 5.3, this results in the following replay sequences, depending on the value of *precrash_txg*.

<i>precrash_txg</i>	Sequence
3	A B C D E F G H I
4	C D E F G H I
5	D F H I
6	F I
7	-

Figure 5.4: Replay sequences for the log depicted in Figure 5.3, by *precrash_txg*.

We use a total order so that we can precisely encode replay progress by storing the last-replayed ($E.gen, E.gsid$). This is necessary for crash consistency because log entry replay is not idempotent. We revisit this topic in Section 5.8.

Note that the definition of the replay order does not account for overflows of *gen* or *gsid* — entries written after the overflow would be incorrectly ordered as smaller than entries written before the overflow. Overflows could be handled in software, e.g., by temporarily destroying the log of the dataset and creating a new one with a fresh log GUID. However, our implementation has no such provisions because even with the (absurdly) conservative of 1 ns write time per log entry, the first overflow event would only occur in 584 years if *gen* starts at 1.

5.6 HDL Replay: Dependency Tracking

Replay is complicated by the fact the entries that were stored in the PRB can be lost. For example, bitflips in PMEM might corrupt the entry’s log GUID or PRB-internal metadata. If an entry E_m with $E_m.txg > precrash_txg$ is missing, any entry E_d that logically depends on E_m ($E_m < E_d$) and is for an unsynced txg ($E_d.txg > precrash_txg$) must no longer be replayed. However, all entries E_p that do not depend on E_m ($E_p \leq E_m$) and need replay ($E_p.txg > precrash_txg$) should still be replayed.

We detect missing entries through a set of counters that we store in the metadata of each entry. For an entry E , the counter $E.C_{ctxg_i}$

$$E.C_{ctxg_i} := \#\{D : D.gen < E.gen \wedge D.txg = ctxg_i\}$$

stores the number of entries that were written to the log since its inception, for a DMU transaction with txg $ctxg_i$, until generation $E.gen$. During replay, we first construct the replay sequence (example in the previous section, Figure 5.4). Then, we re-compute and compare the counters for each entry in the replay

sequence. If an entry E_m has been lost, the counters of any dependent entry E_d (their $E_d.gen > E_m.gen$) will not match, causing replay to stop with E_d as a *witness*. Missing entries in the last generation (the “tail” of the log) cannot be detected with this scheme.

The per-txg scoping of the counters is critical to accomodate garbage collection. Suppose we only used a single sequence counter for all log entries of a HDL. After a txg t_{syncd} has been synced, PRB garbage-collects all entries \mathcal{S}_{gc} that are obsolete:

$$\mathcal{S}_{gc} := \{E : E.txg \leq t_{syncd}\}$$

These entries no longer show up when the claiming or replay procedures scan the PRB for entries with the HDL’s log GUID. If we used a single sequence counter to check for missing entries, we would be unable to discern garbage-collected entries from missing entries.

It is sufficient include only those counters $E.C_{txg_i}$ in the entry metadata whose transaction groups txg_i had not yet been synced at the time that $E.gen$ started. The reason is that a) the replay sequence only contains entries in unsynced txgs and b) E only depends on entries D with $D.gen < E.gen$. Since there are only three possible unsynced txgs at any time (*open*, *quiescing*, *syncing*), we only need to store three (txg_i, C_{txg_i}) tuples per log entry. In fact, since txgs never skip a number, we only store the value of txg_{open} and recompute

$$\begin{aligned} txg_{quiescing} &= txg_{open} - 1 \\ txg_{syncing} &= txg_{open} - 2 \end{aligned}$$

compact encoding is just an idea, not in impl yet

We use the same algorithm to compute the counters on both the write and recovery path. This works because the write path must use monotonically increasing generation numbers, and the entries in the replay sequence are sorted in that order. The counters are stored in a table called *live table*. It has four rows $R_i := (txg, ctr), i \in \{0, 1, 2, 3\}$. When an entry E is written or visited, it modifies the counter in the row with index $I_T := T \bmod 4$ where $T := E.txg$. We distinguish the following conditions:

- If $R_{I_T}.txg = T$, we simply increment $R_{I_T}.ctr$ and are done.
- If $T > R_{I_T}.txg$, we can infer that txg $R_{I_T}.txg$ must have been synced out because there are only three unsynced txgs at any given time but four array entries that are reused in a circular manner, courtesy of indexing by $T \bmod 4$. In that case, we can discard the state in the row and reuse it for T by setting $R_{I_T} \leftarrow (T, 1)$.

- Conversely, if $T < R_{IT}$, we can infer that the entry we are writing is for an already synced txg and hence obsolete. In that case, we do not change the *live table* and turn the entry write operation into a no-op.

Note that the use of 4-ary arrays allows the use of *bitwise-and* for indexing, which is a common ZFS idiom (`table[txg&3]`).

To detect the start of a new generation and support crash-consistent replay, we maintain a separate variable $E_{last} := (gen, gsid)$. When an entry E is written or visisted, we compare it to E_{last} using the following rules:

- $E < E_{last}$ This case is prohibited because the API contract for log writers is that generation numbers must be monotonically increasing. We consider the log corrupted if such an entry E is found.
- $E.gen = E_{last}.gen$ The writer did not start a new generation and we leave the *last table* unmodified.
- $E.gen > E_{last}.gen$ We create a copy of the *live table*. We refer to this snapshot as *last table*.

Regardless of whether a new generation was started or not, we always update $E_{last} \leftarrow E$ and increment the counter in the *live table* as described in the previous paragraph after evaluation the rules above.

The counters $E.C_{txg_i}$ have the same value for all entries in generation $E.gen$ because, by defintion, they only count entries written up to but not including $E.gen$. Hence we compute them from the *last table* once and cache the results until the next generation is started:

1. Determine row index $i_{max} := \max_{i \in \{0,1,2,3\}} R_i.txg$ where $R_i \in$ last table.
2. Invariant: $R_{i_{max}}.txg$ is the highest potentially unsynced txg in generations $< open_gen$. We make the most conservative assumption that $R_{i_{max}}.txg$ was the *open txg* in generation $open_gen - 1$.
3. Find the counters for *quiescing* and *syncing txg*. We scan the *last table* twice to find counters for transaction groups $R_{i_{max}}.txg - 1$ and $R_{i_{max}}.txg - 2$. If the *last table* does not contain those counters, we use a value of zero.

Figure 5.5 contains a graphical illustration of the algorithm described above. Figure 5.6 shows how we encode the counters in the entry header. We provide an extensive example in Figure 5.7.

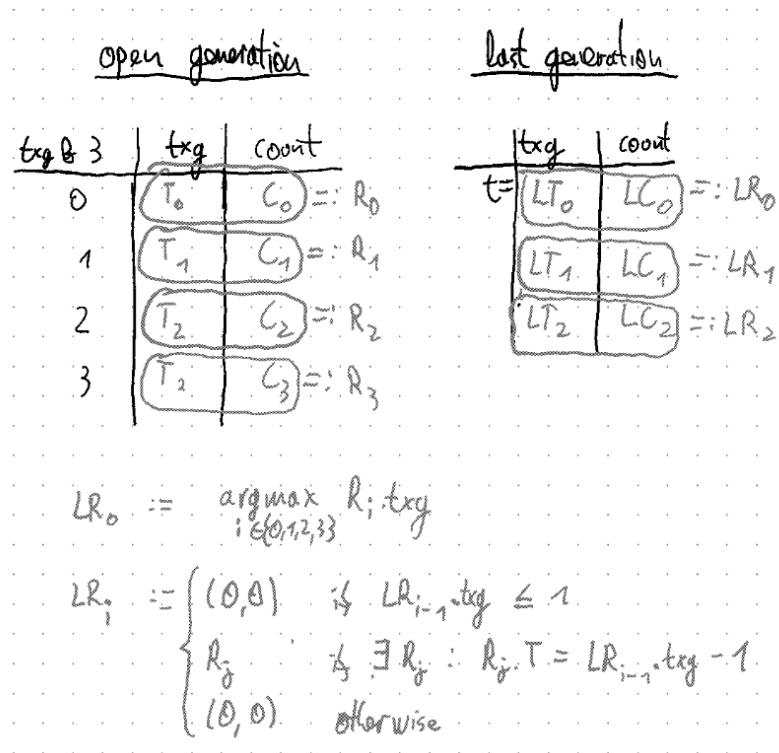


Figure 5.5: Visualization of the *live* and *last table* structures, and a declarative definition of how the *last table* is computed.

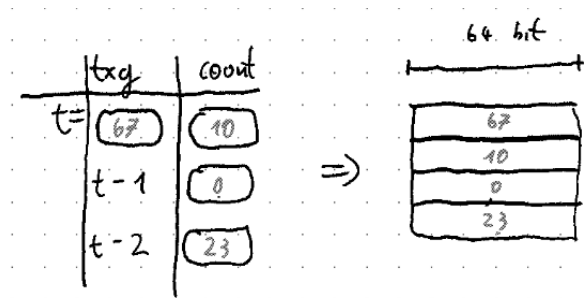


Figure 5.6: Space-efficient encoding of the counters.

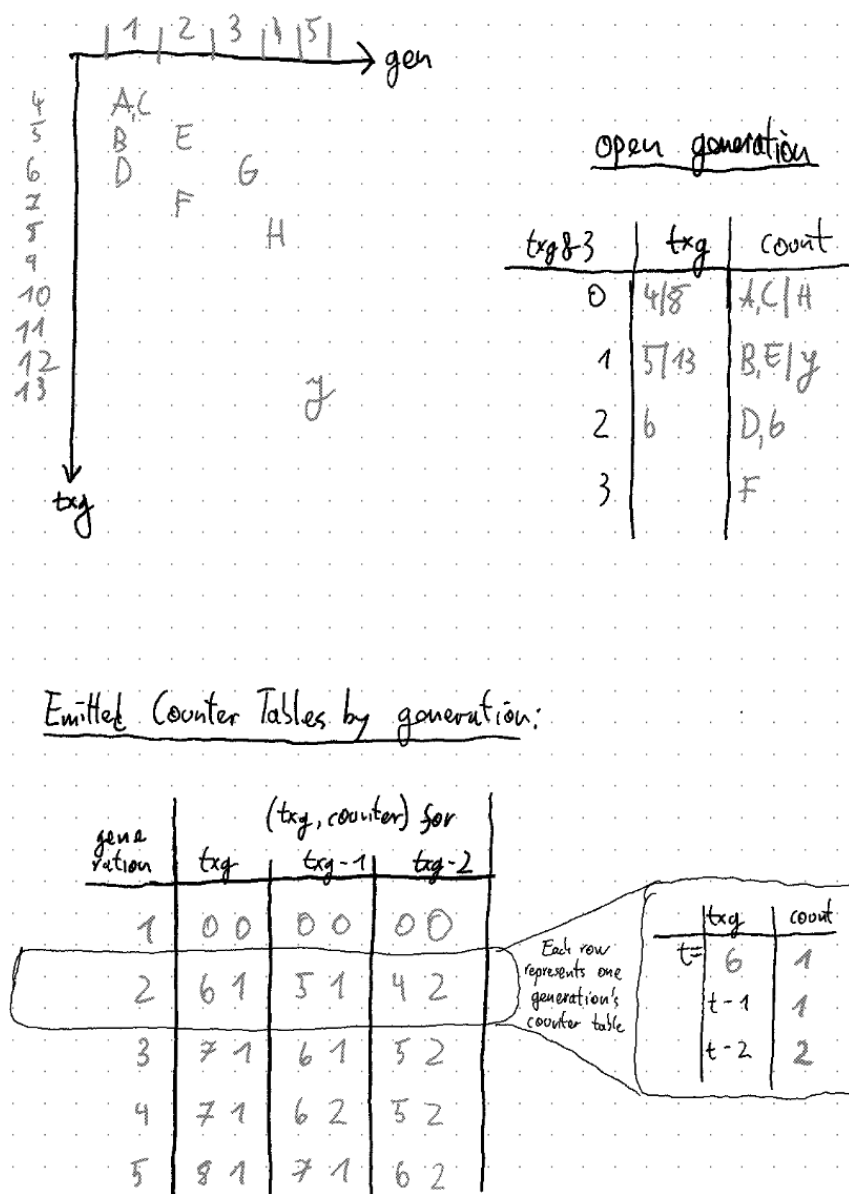


Figure 5.7: Example for the computation of the counters that are stored in the entry. The following entries are written to the log: A,B,C,D in generation 1; E,F in generation 2; G in generation 3; H in generation 4; J in generation 5. Note that their txgs differ sometimes but stay in a corridor of at most three txgs. The table at the bottom shows the counter values in the entry headers: each row represents the counter table that is stored in the entry headers of one generation. The table at the upper right shows the algorithm's table's content over time. Instead of counters, we show the entries that would cause the specific counter to be incremented. We visualize row-reuse by separating new row content with a "|" symbol in each cell of the reused row.

5.7 HDL: Model For Data Corruption

The use of counters to validate that no log entries are missing relies on the assumption that random data corruption is incapable of “forging” new log entries. If this were possible, a forged entry could compensate a lost entry in the counter table. The loss of the entry would go unnoticed and the replay of the forged entry would likely corrupt the dataset.

We believe that it is sufficiently unlikely for random data corruption to accidentally forge an entry. A forged entry would have to fulfill the following requirements to affect counter validation for a given log L :

- It must have been correctly added to PRB’s data structures such that a scan of PRB will find it during recovery. (We address the PRB’s data structure in Section 5.10.)
- Its *log GUID* value must match log L ’s log GUID.
- Its $(gen, gsid)$ must not collide with the original entries. Such a collision would be noticed when constructing the replay sequence. (Our implementation uses a b-tree that identifies and sorts nodes by $(gen, gsid)$).
- Its *txg* value must be in the correct corridor of unsynced txgs and generation. If the txg is too old or too young, the forged entry would be noticed by the dependency tracking code when re-computing the counters (see previous section).

Apart from random data corruption, we have considered the following scenarios and concluded that they are out of scope for ZIL-PMEM:

- Implementation errors on the write path.
- Firmware bugs, e.g., in the Optane DIMM’s wear-leveling layers, that could make old entries re-appear.
- Deliberate injection of log entries by a malicious party with write access to PMEM. This needs to be part of the threat model if ZIL-PMEM is extended to support OpenZFS Native Encryption. Forged entries should be trivially detectable by cryptographically authenticating the entry header, e.g. using the AEAD ciphers that are already in use for the encryption feature.
- Accidental reuse of *log GUIDs*. If two HDLs use the same log GUID to write their entries, the HDLs will each claim their and the other HDL’s entries. This scenario is very unlikely because we generate log GUIDs as 128 bit random numbers and check for collisions with other HDLs.

5.8 HDL: Replay Crash-Consistency

The PRB and its HDLs must be able to tolerate a system crash at any time in their lifecycle. With regards to recovery, this means that claiming must be restartable and replay resumable across crashes. Further, the restarted or resumed recovery procedure must be able to handle online and offline loss of entries due to data corruption.

Crash-consistency for claiming is trivial because it runs during pool import before txg sync starts. Thus, all changes made by claiming are accumulated in DRAM and atomically persisted by the first transaction group that is synced after pool import. From the perspective of HDLs that are in state *logging*, this first transaction group is $precrash_txg + 1$. If the system crashes before that txg is synced, a subsequent import attempt restarts from the same state as the previous one. HDLs that are already in state *replaying* make no modifications during claiming — they only build up DRAM state, i.e., the set of entries held back for replay.

Replay is complicated by two additional aspects:

- Replay of an individual log entry is not idempotent.
- Replay is spread across several transaction groups and hence not atomic. (See this chapter's Section 5.1 on OpenZFS background.)

Our solution is to externalize all state of the replay algorithm, including the counters used for dependency tracking, to a structure called *replay state*. Whenever we replay an entry, we checkpoint *replay state*. We store the latest checkpoint in the ZIL header every time we replay an entry. If the system crashes, we recover *replay state* from the checkpoint in the ZIL header.

The following variables are stored in *replay state*:

(gen_{last} , $gsid_{last}$) The $(E.gen, E.gsid)$ of the last entry E that was replayed.

Live Table The *live table* at the time the entry was replayed, i.e., R_i with $i \in 0, 1, 2, 3$.

Last Table The *last table* at the time the entry was replayed.

E_{seal} An artificial entry header that is used to detect lost entries at the end of the log structure.

The crash-safe replay procedure performs the following steps:

1. Claiming reads the log GUID from the ZIL header, then scans the PRB and holds back the log's entries from garbage collection.

- If the HDL is in state *logging*, claiming transitions the ZIL header to *replaying* and initializes the *replay state* checkpoint in the ZIL header.
 - $(gen_{last}, gsid_{last})$ is $(0, 0)$.
 - *live table* and *last table* are zeroed out.
 - E_{seal} is determined as follows.
 - (a) Construct a temporary replay sequence $S_{tmp} = \{E_s : s \in 1, \dots, n\}$.
 - (b) In DRAM, append an artificial entry E_{n+1} to S_{tmp} . $E_{n+1}.gen = E_n.gen + 1$, $E_{n+1}.gsid = 1$, and $E_{n+1}.txg = E_n.txg$.
 - (c) Iterate over S_{tmp} and compute the counter tables. Validate the tables for $E_{1,\dots,n}$. When arriving at E_{n+1} , assign the expected table to E_{seal} .
 - If the HDL is already in state *replaying*, the ZIL header is not modified.
2. Replay recovers its *replay state* from the checkpoint C in the ZIL header.
 3. Replay constructs a replay sequence S from the held back entries.
 4. While replaying S , the following happens for each entry $E_s \in S$:
 - (a) Skip E_s if it has already been replayed, i.e., $(E_s.gen, E_s.gsid) \leq (C.E_{last}.gen, C.E_{last}.gsid)$. Otherwise:
 - (b) Create a backup checkpoint C_b of *replay state*.
 - (c) Perform dependency tracking as described in the previous section.
 - (d) Create a checkpoint C_{E_s} .
 - (e) In one DMU transaction, replay E_s and store C_{E_s} in the ZIL header.
 - (f) If replay fails or the DMU transaction fails, abort the transaction and roll back the in-DRAM version of *replay state* to C_b .

Note that steps 1 and 3 construct a *new* replay sequence every pool import. This allows handling of offline data corruption in PRB. Assume that we lose an entry E_m while the system is offline. If $E_m.txg \leq precrash_txg$, the loss of the entry is unnoticed because E_m is by definition not part of the replay sequence. If E_m is skipped by step 4a, the loss of the entry might be worth reporting but does not affect replay because it has already been replayed. Otherwise, dependency tracking will eventually detect that E_m is missing if, and only if, any other entry depends on it. E_{seal} is an artificial dependency to ensure that lost entries in the last generation are detected.

this is not yet implemented

We handle online data corruption as follows. Claiming and replay only operate on DRAM-buffered copies of the entry metadata called *replay node*. When replaying an entry, the replay callback is forced to use a special function that attempts to read the entry from PRB and buffers it in DRAM. Before allowing access to the entry, it ensures that the entry metadata still matches the data in the *replay node*. If either the read from PRB fails (e.g., due to checksum errors)

or if the metadata does not match, we know that the entry has been corrupted since the PRB scan.

Note that the last generation $E_{seal}.gen - 1$ of the log structure is not protected against re-appearing log entries. Assume that during initial claiming, entries (42, 1) and (42, 3) in generation 42 are observed as the last entries of the replay sequence. Then $(E_{seal}.gen, E_{seal}.gsid)$ is (43, 1) and the counter table of E_{seal} only accounts for these two entries, and all previous generation's entries. Now assume that another entry (42, 2) had been written but was temporarily invisible during initial claiming, e.g., due to temporary data corruption. If for some reason (42, 2) becomes visible again, and we start replaying, then (42, 2) will be replayed even though it was not visible during claiming. And if (42, 2) were to replace (42, 1) or (42, 3), we would not even notice the counter mismatch when arriving at E_{seal} . This behavior seems desirable in the sense that replay recovers as much committed data as possible. However, in combination with WR_INDIRECT ZIL log records, it can lead to pool corruption. ZIL-PMEM does not currently support WR_INDIRECT records due to this problem as well as another WR_INDIRECT-related issue that affects ZIL-LWB. Our proposed solution in the OpenZFS issue tracker could be used by ZIL-PMEM to properly detect changes to the last generation after claiming, and to support WR_INDIRECT if desired.

link github
issue on leak-
ing allocated
blocks

5.9 PRB: DRAM Data Structure

The central requirements for PRB are

- persistence of log entries in PMEM,
- parallel insertion of log entries from multiple threads and HDLs,
- garbage collection of obsolete log entries,
- measures to detect data corruption.

Our design is built around the *chunk* abstraction. A chunk is a contiguous segment of PMEM that acts as an insert-only container for log entries. PRB's role is to mediate access between HDLs and chunks. When a thread writes an entry to the conceptual *set of entries* that is the PRB, it actually inserts the entry into one of the PRB's chunks. Once a chunk is full, PRB lets it sit unmodified until all entries in it are obsolete. Garbage collection then removes all entries in those obsolete-only chunks and makes them available for reuse by writers. The PRB scanning done by claiming and replay is implemented as an iteration over the entries in all chunks of the PRB. If a chunk contains at least one candidate entry for replay, the claiming HDL holds the chunk back from the write path and/or

garbage collection. Note that this design fully decouples storage location (chunk) from log structure (encoded in entry metadata, see previous sections).

Each chunk is represented by a DRAM object that holds the following values:

PMEM location The chunk’s start and end address in PMEM. These values do not change for the lifetime of the chunk.

Claim refcount The number of HDLs that claimed at least one entry in the chunk for replay. Chunks with non-zero claim refcount are exempt from garbage collection and the write path.

Max txg The maximum transaction group of the entries that were written to this chunk. This value is used by garbage collection to determine when all entries in the chunk are obsolete.

Write position The PMEM address where the next entry will be written. We discuss the PMEM layout of the chunk in more detail in Section ??.

We move chunks between the following data structures to keep track of their current role within PRB.

Commit Slots PRB maintains a roster of *commit slots*, each of which has an associated *open chunk*. When a thread writes an entry through a HDL, the thread first *acquires* a commit slot to gain temporary exclusive access to its *open chunk*. Then, it writes the entry to the open chunk and immediately releases the slot so that other threads can acquire it.

Free List The *free list* holds empty chunks. Their *claim refcount* and *max txg* is zero and their write position is reset to the chunk’s start address. A log writer that needs a new chunk for its commit slot puts the commit slot’s current open chunk on the *full list* (see below) and gets a new chunk from the *free list*.

Full Lists The *full lists* contains chunks that wait for garbage collection. PRB maintains one full list for each unsynced txg T_i . The *full list* for T_i contains exactly those (full) chunks whose *max txg* is T_i . After T_i has been synced, the chunks on the corresponding *full list* are emptied and moved to the *free list*. Note that it is sufficient to maintain three full lists at any given time — one for each of the three unsynced transaction groups *open*, *quiescing*, and *syncing*. The three lists can be represented as an array of three lists that is indexed by $T_i \bmod 3$. Our implementation follows the common ZFS idiom of using four lists that are indexed by $\tau_i \& 3$ where $\&$ is the *bitwise and* operation.

Wait Replay List During the claiming phase, all of the PRB’s chunks are on the *wait replay list*. The HDL’s scan the chunks on this list during claiming. If a chunk contains at least one entry that needs to be held back for replay, the HDL increments the chunk’s *claim refcount*. Once claiming is done

and *txg sync* triggers the first garbage collection cycle, any chunk on this list that has a zero refcount is emptied and moved to the *free list*. The list is also scanned for zero refcounts during all future garbage collection cycles to garbage-collect chunks that are no longer held by HDLs because they finished replay. Note that it is critical for replay crash consistency to defer garbage collection until *after* the last txg of replay of the last HDL that held the chunk has synced. If we garbage-collected entries before the last holding HDL's ZIL header has transitioned to state *nozil* on disk and crashed inbetween, the garbage-collected entries might be missing when resuming replay.

Figure 5.8 visualizes transitions of the chunks over their lifetime. Figure 5.9 provides a corresponding example.

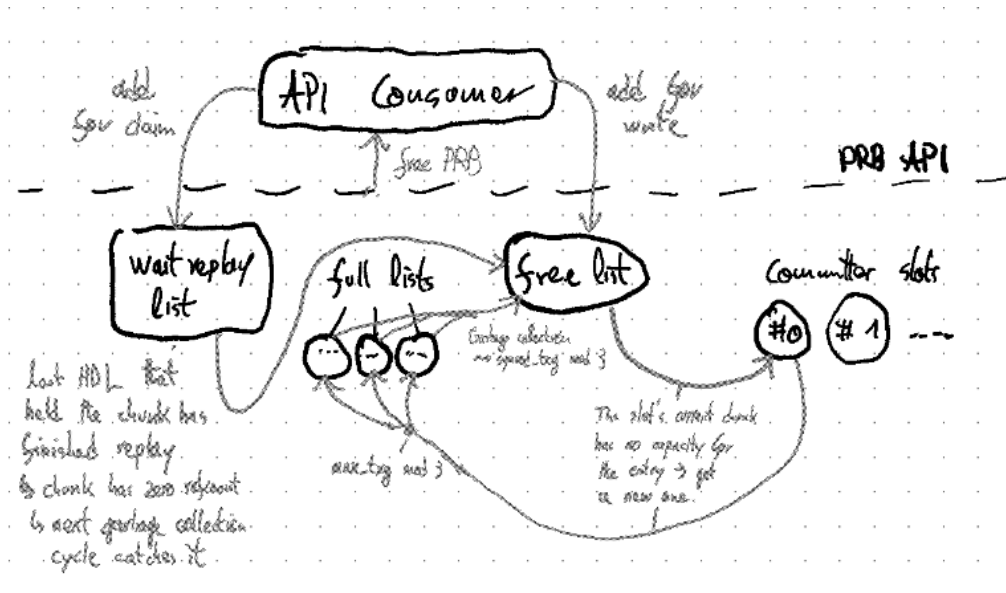


Figure 5.8: The different owners of a chunk and the events that cause ownership transitions.

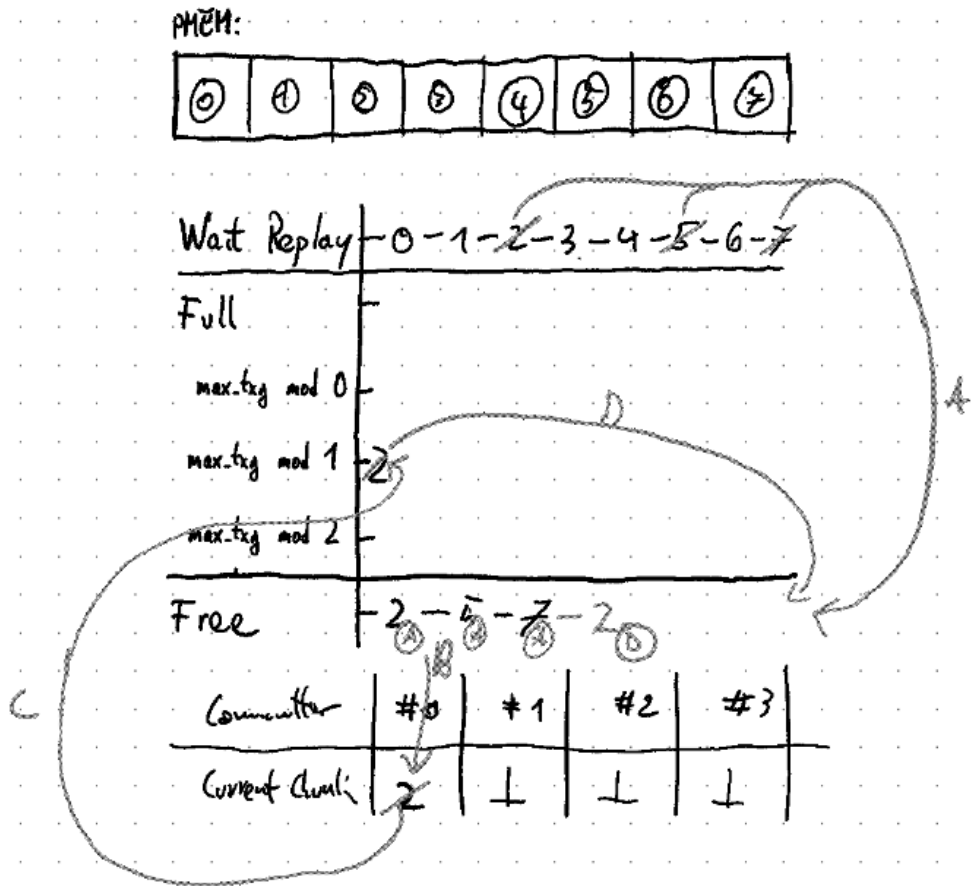


Figure 5.9: Example for chunk ownership transitions. The PRB is constructed with all eight chunks added as `_for_claim`. Claiming (A) determines that chunks 0, 1, 3, 4, and 6 need to remain on the *wait replay list* because they contain entries for HDL logs that need replay. Chunks 2, 5, and 7 only contain obsolete entries and move to the *free list*. We do not replay any of the logs. However, a log writer starts writing a new log in B. It finds that its commit slot (#0) has no active chunk. Thus, it and moves the first chunk on the *free list* (i.e., chunk 2) to the commit slot. After writing the entry to chunk 2 the log writer releases the commit slot. Another log writer acquires commit slot #0 and writes entries to it. Chunk 2's capacity is insufficient to hold the last entry (C). The thread places chunk 2 on the correct *full list* for chunk 2's *max txg* and finds a new chunk on the free list for the commit slot (not shown). Eventually *txg sync* triggers garbage collection for the *txg* that is chunk 2's *max txg* which resets chunk 2's in-PMEM and in-DRAM sequence and subsequently places it back onto the free list.

5.10 PRB: PMEM Data Structure

As explained in the previous section, PRB is built on top of contiguous slices of PMEM which we call *chunks*. However, PRB only consumes the chunks, it does not create them. Chunk allocation is the responsibility of the PRB consumer, i.e., ZIL-PMEM. This design improves modularity and testability because it decouples the following two concerns:

Resource Aquisition The PRB consumer is responsible for integrating PMEM SLOG vdev into the zpool, discovering its memory mapping, and partitioning the PMEM space.

PMEM Data Structure PRB only implements the persistent data structure that stores entries in the chunks (next section).

ref

With regards to persistent data structure, this separation of concerns relieves PRB from the need to define structures to track the partitioning of PMEM space into chunks. It relies on the consumer to provide the PRB with the same set of chunks every time the PRB is constructed.

Entries are variable-length records that are stored within the chunk as a contiguous sequence. Each entry is represented as a fixed-length 256 byte sized header and a variable-length body. The first entry starts at the chunk's start address which must be aligned to 256 bytes. The space after each entry is zero-padded to the next multiple of 256 bytes. The next entry starts after the padding. The sequence is terminated either explicitly by an invalid entry header or implicitly if the last entry has filled the chunk completely. Figure 5.10 provides an example chunk layout. Note that the ordering of entries within the chunk is has no semantic value as the logical view on the chunk is that of a set of entries.

The entry body is a verbatim copy of an opaque byte slice provided by the log writer. (In ZIL-PMEM, this byte slice is the log record structure that is shared among all ZIL kinds.) The entry header's contents are managed by the PRB and HDL. Its contents are as follows:

language

HDL-scoped metadata The metadata required for attribution of a log entry to a HDL and subsequent replay.

- Log GUID
- Generation
- Generation-Scoped ID
- Encoded Counters for dependency tracking.

Body Length We store the exact body length in bytes. The zero padding in the chunk sequence is not considered part of the entry itself.

fix that?

Body Checksum Fletcher checksum of the body data.

Header Checksum Fletcher checksum of the header to ensure data integrity of the metadata. If the header checksum is corrupted then none of the other header fields can be trusted.

Zero Padding The unused bytes in the header have the defined value of zero. Their value is part of the header checksum.

Chunks must be sufficiently large to hold at least one entry because there is no mechanism to split an entry across multiple chunks. The smallest chunk's size determines the maximum entry size that can be written to it. However, chunk sizes much larger than a single average entry are advisable for multicore-scalability, as we will elaborate on in Section 5.13.

Basic In-PMEM Format

The PMEM DIMM is partitioned into **chunks**.
A chunk has a 256B-aligned base address.
Chunks contain a contiguous sequence of **entries**.
The first entry starts at offset zero.

An **entry** consists of

- a fixed-length **header**
- a variable-length **body**
- padding to the next 256B multiple.

The header contains the **body's length** and **checksum**.
A separate checksum protects the header fields.

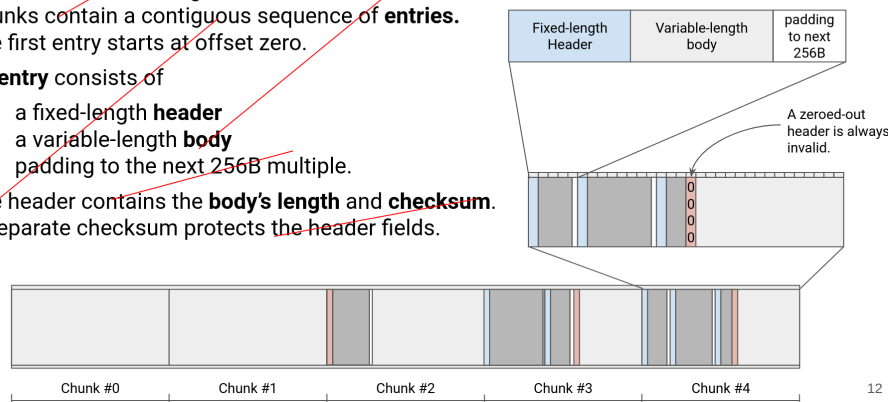


Figure 5.10: Example chunk layout. Note that although we partition the PMEM space very regularly in this example, the PRB consumer is free to use variable-sized chunks.

5.11 PRB: Chunk Traversal

HDLs scan the PRB during claiming to put their holds on chunks that contain replayable log entries. During replay, the HDLs scan their held chunks again to construct the replay sequence. The building block for both procedures is the iteration over the entries of a chunk. We call this iteration process *chunk traversal* and present the algorithm in this section.

As a reminder, the data layout within the chunk is as follows:

- The first entry header starts at chunk offset zero. Its size is fixed (256 bytes).

- The variable-length body starts immediately after the header. Its length is stored in the header.
- The body is followed by zero padding to the next 256 byte multiple.
- The next entry's header starts there.
- An invalid header marks the end of the chain. A header is invalid if its header checksum is invalid, or if the log GUID is invalid. The (128 bit) log GUID is invalid if at least the upper or lower 64 bits are zero.
- Entries are never split between chunks.

In the absence of data corruption in PMEM, chunk traversal visits each valid entry in the sequence and stops at its end. For the case where PMEM is corrupted, we distinguish the following conditions:

word?

Machine check exception (MCE) If the PMEM hardware detects uncorrectable data corruption, it raises an MCE. The Linux kernel provides the an API (`memcpy_mcsafe`) which provides a `memcpy`-compatible interface but converts MCEs into error return values. We always use this API to buffer PMEM contents in DRAM before accessing them. The error handling for MCE errors is the same as for invalid checksums in entry header and body, which we describe below.

Header: Detected Data Corruption If the header checksum validation fails, the header was either corrupted by bitflips or similar phenomena or never completely written. The latter case is critical for crash consistency on the write path as we will elaborate on in Section 5.13.3. Regardless of the cause, an invalid header's values cannot be trusted, and thus the traversal stops.

Header: Undetected Data Corruption If the header checksum does not detect data corruption in the header, the behavior is implementation defined. However, it is guaranteed that memory accesses are constrained to the entry's chunk's bounds.

Data corruption in the body The traversal algorithm does not read the body but instead returns a closure that can be invoked for this purpose. The closure reads the body into DRAM using `memcpy_mcsafe`, then validates the body checksum stored in the header. Validation failures are returned as an error. The caller can decide whether they want to iterate further or propagate the error up the call stack. Note that in our C implementation, the closure is replaced by the opaque struct `zilpmem_replay_node_t` and function `zilpmem_prb_replay_read_replay_node`.

Data corruption in the padding We require that the padding in the space that follows the body to the next 256 byte multiple consists only of zero. We validate this property in the closure that reads the body. Validation failure

results in a distinguished error being returned from the closure. Such an error is indicative of data corruption or an incorrect implementation on the write path. The caller of the traversal algorithm should surface it to the administrator, but may choose to proceed.

The following listing provides the pseudo code for chunk traversal.

Algorithm: `iter_chunk()` - Iterate over the entries in a chunk.

Inputs:

`ch_base` chunk's base address
`ch_len` chunk's length

Procedure:

```

assert ch_base % 256 == 0
assert sizeof(entry_header_t) == 256
assert ch_len >= 256
assert ch_len % 256 == 0

e := ch_base
while (e < ch_base + ch_len) {
    entry_header_t eh;

    // read header with a function that handles MCEs
    memcpy_mcsafe(&eh, e, sizeof(eh));

    // invalid hdr checksum or id terminate the sequence
    validate_header_checksum(eh)?;
    if eh.log_guid's upper or lower 64 bits are zero {
        return None;
    }

    body_ptr := e + sizeof(eh);

    if body_ptr + eh.body_len > ch_base + ch_len {
        return Err(
            "entry out of chunk boundary:
             a) writer implementation error
             b) undetected header corruption
            ");
    }
    e += sizeof(eh) + body_len;
    e = roundup_to_next_256byte_multiple(e);

    padding_ptr := body_ptr + eh.body_len
    padding_len := e - padding_ptr

    read_body := |buffer: *u8| {
        memcpy_mcsafe(buffer, body_ptr, eh.body_len)?;
        validate_body_checksum(eh, buffer, eh.body_len)?;
        pmem_is_zero(padding_ptr, padding_len)?;
        return Ok();
    };
    yield Some((eh, read_body))
}

```

Example:

```

    for (entry_header, read_body) in iter_chunk(...) {
        if entry_header.txg <= precrash_txg {
            continue;
        }
        buf := buffer of size entry_header.body_len
        read_body(buf)?;
        ...
    }

```

5.12 PRB: Garbage Collection

We already covered the essence of garbage collection in Section 5.9: if the capacity of an acquired commit slot is insufficient, the log writer puts it on the *full list* for the chunks' *max txg* and gets a new chunk from the *free list*. After *txg sync* has written out that *max txg*, it triggers garbage collection, which removes all entries in the chunks on the *full list* for *max txg*.

We remove all entries in a chunk in $O(1)$ time by invalidating the log GUID in the first entry header in the chunk (offset zero). As indicated by the This is sufficient to prevent chunk traversal in the future. The pseudo-code to invalidate a chunk is as follows:

```

Input:
    ch_base      chunk base address
    ch_len       chunk length
Steps:
    assert ch_len >= 256

    assert type of entry_header_t::log_guid is uint64_t[2]
    log_guid_addr := ch_base + offsetof(entry_header_t, log_guid)
    assert log_guid_addr is 64-bit aligned

    atomic 64-bit store of 0 to log_guid_addr
    atomic 64-bit store of 0 to log_guid_addr + 8

    flush modified cache lines
    sfence

```

5.13 The Write Path

A thread that writes an entry to a HDL needs to perform the following tasks:

- Find a target chunk using the *commit slot* mechanism. (Section 5.13.1)
- Determine the HDL-scoped metadata, i.e., log GUID, txg, gen, gsid, and counters.
- Compute header and body checksums.
- Insert the entry into the target chunk in a crash-consistent manner.

Our goal is a design with low write latency, good multicore-scalability, and CPU efficiency. We identify the following factors as particularly relevant:

Checksumming Checksumming adds latency but does not concern multicore scalability since no coordination is required between writers. The implementation should use ZFS's optimized implementations of the Fletcher checksum.

Dependency Tracking Counters We must update the dependency-tracking counters on every entry write operation. Since the counters are HDL-scoped, this only presents a scalability concern if a single HDL is written from multiple threads. Parallel writes to the same HDL do not happen in ZIL-PMEM proper but are relevant for our ZVOL-specific ITXG bypass (Section 6.3.5).

Commit Slot Aquisition & Chunk Replacement The PRB's commit slots and chunk lists are shared among all HDLs. All threads that write to any HDL of the PRB compete for these resources, making it a multi-core scalability challenge.

Optane Characteristics We develop and evaluate ZIL-PMEM for/on Intel Optane DC Persistent Memory. The performance characteristics of Optane DIMMs are significantly different from regular DRAM. Yang et al. have established that the Optane PMEM hardware is organized in units of 256 bytes. For example, the access granularity and kind of store and cache flush instruction have significant impact on the achievable write bandwidth. For size of log entries written by ZIL-PMEM, the use of AVX-512 non-temporal store instructions is recommend for highest possible performance. [Yan+20].

PMEM Bandwidth Limits & Multicore Scalability It is inherent to the programming model for persistent memory that wait time for PMEM I/O is spent on-CPU. For example, instructions that architecturally depend on a preceding store+cacheflush+sfence to PMEM will stall until the flushed cacheline reaches the power-fail protected domain of the CPU [Sca20]. This is problematic from a CPU utilization perspective: if multiple threads attempt to write at higher bandwidth than PMEM can sustain, they still appear busy towards the OS thread scheduler and waste CPU time that could be used more productively by other threads in the system. Yang et al. have shown that a single Optane DIMM's write bandwidth can be exhausted by one CPU core at 2 GB/s. Write bandwidth decreases to 1 GB/s at ten or more concurrently writing CPU cores [Yan+20]. Whereas excessive on-CPU waiting might be the right trade-off in certain userspace applications of PMEM, a kernel file system such as ZFS cannot make assumptions about the system's overall CPU priorities. We expect that PMEM write

bandwidth can be exhausted in real-world use cases for ZIL-PMEM. Our design must therefore find a way to limit concurrent access to PMEM and shift PMEM wait time off the CPU.

5.13.1 Commit Slots

Commit slots are our abstraction to enable multiple threads to write entries concurrently. The goal is to grant up to `ncommitters` parallel writers temporary exclusive access to a chunk into which they can write their log entry. For this purpose, a thread that wants to write an entry *aquires a commit slot* $S \in 0, 1, \dots, ncommitters - 1$. Each thread that is simultaneously committing gets a different commit slot. If no commit slot is available, the function to aquire the commit slot blocks. Associated with each commit slot is a PMEM chunk which we refer to as *open chunk*. A thread that aquires a commit slot is allows to write its entry to the slot's *open chunk*.

The limitation to `ncommitter` parallel writers is desirable to avoid excessive on-CPU waiting on PMEM. Let us make the simplifying assumption of a fixed maximum write bandwidth of $B_{max}[byte/s]$ to PMEM before latency increases dramatically due to queuing. Then an equal distribution of that bandwidth yields $\frac{B_{max}}{ncommitters}[byte/s]$ of write bandwidth per writer. Assuming that writers do not exceed this limit, we can derive latency guarantees for writing entries, dependent on entry size.

We implement commit slots using a semaphore initialized to `ncommitters` and a bitmask with `ncommitters` bits. The thread that aquires a commit slot first enters the semaphore and then finds and flips a zero bit in the bitmask. The zero bit's index is the commit slot number. We use opportunistic spinning to find and flip the bit. The following pseudo-code explains the aquisition and release procedures.

```

Procedure For Commit Slot Aquisition:
Input:
    sem        semaphore
    bm         pointer to PRB-wide bitmask with ncommitters bits
Output:
    The commit slot number.
Steps
    Enter semaphore.

    my_bm      <-  atomic_load(bm, SeqCst)
'retry:
    idx        <-  find first set bit index in (~my_bm)
    if idx == 0:
        panic: idx=0 indicates there is no free bit,
                but the semaphore guarantees that
    idx -= 1
    if idx >= ncommitters:
        panic: semaphore guarantees that there are
                free commit slots

```

```

my_bm = my_bm | (1<<idx)
if compare_and_swap(bm, &my_bm, SeqCst):
    // we won the race => return the commit slots
    return idx
else:
    // we lost the race with another committer
    // => retry
    // (my_bm contains the actual value of bm)
    goto retry
---
```

Procedure For Commit Slot Release:

Input:
 cslot The commit slot number returned on aquisition.

Steps:
 Atomic bitwise and of bm with ~(1<<idx)
 Exit Semaphore

The procedure above is all the PRB-wide coordination that is required for writing a log entry, iff the aquired slot's *open chunk*'s capacity is sufficient. If capacity is insufficient, the writing thread replaces the full *open chunk* with a new one. To accomplish that, it puts the current *open chunk* on the correct *full list* and gets a new *open chunk* from the *free list*. Access to the PRB's lists is protected by a PRB-wide mutex. The potential contenders for this mutex are up to *ncommitters* writer threads and the *txg sync* thread that performs garbage collection.

Getting a new chunk from the *free list* fails if the free list is empty. The log writer can choose whether to block and wait or fail the log entry write operation with an error. Block-and-wait is implemented through a condition variable that is signalled by garbage collection for every chunk that it puts back on the *free list*. If the log writer chooses to block and wait during chunk replacement, it must guarantee that it does not prevent *txg sync* from making progress in order to avoid a pool-wide deadlock. In practice, this means that the log writer must not hold a DMU transaction open when writing an entry. Note that this is the case for ZIL-PMEM proper because `zil_commit` is called after the DMU transactions have finished or failed. However, for the ITXG bypass for ZVOLs (Section 6.3.5), we write log entries from within the DMU transaction and thus need to use the non-blocking mode.

Sharing the slots (and thus chunks) among all threads and HDLs in the PRB is of ambiguous value from perspectives other than bandwidth limitation:

Space Efficiency Sharing chunks causes entries to be packed into a small number of chunks, even more so if we implemented some form of *best fit* selection scheme for commit slots. Packing is beneficial for space efficiency

because the *spread* between minimum and maximum txg of the entries in a chunk is small, enabling timely garbage collection.

Blast Radius However, the concentration of entries in a chunk also increases the blast radius of data corruption due to the chunk's physical data structure which we discuss in Section 5.10. In particular, data corruption within the entry metadata of one HDL's entry can render another HDL's entry unreachable during PRB scan if they are stored in the same chunk.

one sentence on limited slog size, that it's not as relevant for PMEM, but maybe for NVDIMM-N?

Cache Efficiency Our acquisition procedure deterministically picks the lowest available commit slot. It is thus very likely that chunks bounce around CPU cores, without taking temporal locality into account. Note that, due to the use of non-temporal store instructions when writing log entries to PMEM, this only impacts the chunk DRAM object. We briefly experimented with per-core commit slots during development and found no noticeable performance difference to the approach presented above.

this last paragraph doesn't feel right here...

5.13.2 HDL-Scoped Metadata

We have already-addressed the algorithm to compute the HDL-scoped entry metadata (Log GUID, txg, gen, gsid, Counters Table) in Section ???. This subsection only addresses multithreading and scalability concerns.

check matches

First of all, HDL-scoped metadata does not pose a scalability concerns if entries are written sequentially. This is the case for ZIL-PMEM proper which uses a mutex to serialize `zil_commit` calls. (Remember from Section 5.1 that the ZIL's shared `itxg` structure defines the sequential *commit list* model.) However, for the ITXG bypass for ZVOLs (Section 6.3.5), the ZVOL dataset's HDL can be written in parallel.

Multiple threads are allowed to write to a single HDL simultaneously iff they do not start a new generation. Threads that start a new generation must wait for all threads that wrote entries to the previous generation to finish writing. The reason is that the new generation's entry logically depends on the previous generation's entry. It would be sufficient to prevent the *function calls* from returning out of dependency order while allowing the new generation's entry to be written in parallel with the old entries. Such a system is used by ZIL-LWB's *commit ITXs* which uses condition variables to wake `zil_committing` threads up when the last LWB that contains one of their commit list's records has been written. However, the *commit ITX* model was not directly applicable to ZIL-PMEM proper. For the ITXG bypass, we use a simple read-write-lock which we elaborate on in Section 6.3.5.

check content

Within the HDL, we use a spinlock to serialize access to the state used for dependency tracking (*live table*, *last table*). We only compute the encoded representation of the *last table* when the first entry is counted for a new generation. Subsequent writes for the same generation only need to `memcpy` the pre-computed encoded representation while holding the spinlock.

5.13.3 Crash-Consistent Insert

After the commit slot has been acquired, a target chunk been selected, and the HDL-scoped metadata determined, we are ready to insert the entry into the chunk. Remember from Section 5.10 that the chunk is a sequence of entries that is terminated by an invalid entry header. To insert the entry into the chunk, we must append the entry to the sequence in a way that is crash-consistent with regards to the traversal algorithm (Section 5.11):

Appending an entry E_{n+1} to a chunk C that contains a sequence of entries E_1, \dots, E_n must be atomic from the perspective of traversal. After a system crash or power failure during the append operation, traversal of C must either find E_1, \dots, E_n or E_1, \dots, E_n, E_{n+1} .

The following algorithm describes the procedure and invariants during the append operation. Figure 5.11 contains a step-by-step visualization.

```

Inputs:
    ch  The chunk into which we append
    e    The entry that we append to ch's sequence.
Procedure:
    Invariant 1: Traversal will stop at ch.pos because
                  either: ch.pos points to PMEM space that
                        is an invalid header
                  or: there is no space left in the chunk

    Phase 1:
        Write e.body to address ch.pos + 256
        Write trailing zero padding
        If there is still space in the chunk at
            address ch.pos + 256 + e.body_len + padding_len:
                Assert that the space is at least 256 bytes.
                Write an invalid follow header (256 zero bytes)
                    to that address.
        Cacheflush + Sfence

    Corrolary 1: Neither body nor follow header
                  will be visited by traversal
                  because traversal still stops
                  at address ch.pos (Invariant 1).

    Phase 2:
        Write e.header (256 bytes) to address ch.pos
        Cacheflush + Sfence

    Corrolary 2: Traversal visits the entry written

```


to <code>ch.pos</code> and stops at the follow header.

Invariant 1 can be proven by induction: if the entry being written is the first entry in the sequence (base case) the chunk was fetched from the *free list* which is defined to only contain chunks that are *empty*. (*Empty* means that the write position (`ch.pos`) is at the chunk's base address and that the PMEM at the write position is an invalid header.) The induction step is that if an entry E_{n+1} is appended to an existing sequence E_n , invariant 1 holds as well. This is true because the space occupied by E_{n+1} contains an invalid follow header written by phase 1 when E_n was appended to the sequence. Stores made in Phase 1 for E_n are guaranteed to be persistent before any store for E_{n+1} happens due to the `cacheflush+sfence` at the end of Phase 1. Note that we do not need to address the case where the chunk is full because we cannot append E_{n+1} to a full chunk.

The first `sfence` in phase 1 is required for correctness iff we do not trust the body checksum to detect partial writes. If we omitted the `sfence` in phase 1, it would be possible that the entry header written in phase 2 reaches PMEM completely but the body written in phase 1 does not. In that case, after a crash, the traversal algorithm would observe a valid header but a body space with undefined content. We would rely fully on the body checksum to detect the partially written body. If the checksum is too weak, the replayer's interpretation of the undefined body contents determines subsequent behavior. Corruption of the dataset is a likely consequence.

The `sfence` (phase 2) is required for correctness because we must conservatively assume that the log writer's subsequent store instructions (in program order) depend on the entry having reached stable storage.

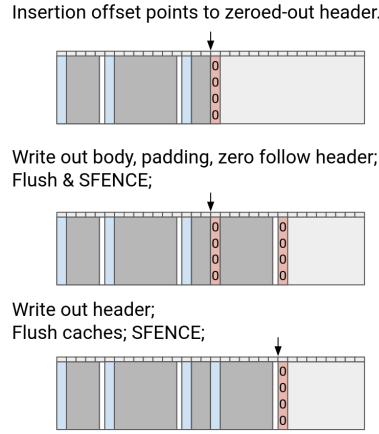


Figure 5.11: The crash-consistent append operation to a PMEM chunk. The traversal algorithm reaches the existing entries at all times and reaches the new entry only after its body and header have been written completely.

The following implementation details of the append operation are relevant for latency and efficient use of CPU time:

- For ZIL-PMEM, the use of `sfence` in phase 1 came with negligible impact on overall latency during development. We assume that this is because the wait time added by the `sfence` is negligible compared to the write time for the entry body. For example, the entry body and zero padding for a 4k sync write is $\text{sizeof}(lr_write_t) + data + padding = 4096 + 192 + 64 = 4352$ bytes large. Assuming 2 GiB/s write bandwidth for a single Optane DIMM [Yan+20], the write time for the entry body is ~ 2 μ s. In contrast, the derived latency for a 256 byte write at that rate is ~ 0.12 μ s. If we use this value as an approximation for the cost of the `sfence`, its latency contribution is only $\frac{0.12}{0.12+2} = 5.6\%$.
- All writes to PMEM happen in multiples of 256 bytes because 256 bytes is Optane's internal write unit size. Writes below this size cause read-modify-write cycles in the hardware and thus cost performance [Yan+20; Zha+21].
- We use AVX-512 non-temporal store instructions (`movnt`) instead of regular stores and cache flushes. Again, this addresses established performance properties of Intel Optane DC Persistent Memory [Yan+20].
- We also use ZFS's optimized implementation of the Fletcher checksum to compute body and header checksums. Whereas the best implementation is chosen by benchmark dynamically at runtime, it is safe to assume that some SIMD ISA extension such as AVX-512 will be used if available.

- OpenZFS cannot rely on the Linux kernel's interface for saving FPU state due to licensing issues []. Unless the FPU is used by a dedicated kernel thread, ZFS must temporarily disable preemption, mask local interrupts, and manually save FPU state. PRB is written directly from the task that calls `zil_commit` and thus incurs this overhead. We refactor ZFS's FPU state management abstraction so that FPU context is only saved once for both checksum computation and writing to PMEM. The time that is spent in this suboptimal state is bounded by the maximum log entry size.

we have no evaluation on impact of interrupt masking, future work

5.14 API Design

We briefly discuss the PRB/HDL API that is used by the `zil_pmem.c` module. In the implementation, most types and functions are prefixed with `zil_pmem_` or `zil_pmem_prb_t` which we omit for brevity.

5.14.1 PRB Setup

```
prb_t* prb_alloc(size_t ncommitters);
void prb_free(prb_t *b, bool free_chunks);

chunk_t* chunk_alloc(uint8_t *pmem_base, size_t len);
void chunk_free(chunk_t *c);

void prb_chunk_initialize_pmem(chunk_t *c);
void prb_add_chunk(prb_t *prb, chunk_t *chunk);
```

The `zpool` import procedure allocates the PRB using the `prb_alloc` function. The returned `prb_t` is owned by the caller which is responsible for freeing it using `prb_free` during pool export.

The PRB consumer allocates the chunk objects using `chunk_alloc`. It adds the chunk to the PRB using `prb_add_chunk`. The PRB assumes ownership of the chunks that are added to it. When the PRB is constructed for the first time, the PRB consumer must call `prb_chunk_initialize_pmem` to reset the PMEM sequence to an empty state.

The `free_chunks` argument to `prb_free` determines whether the PRB should free the chunks objects that were added to it, or whether the ownership moves back to the caller. Note that freeing the chunk object (`chunk_free`) does not alter the chunk's PMEM state.

HDL Setup

```

void zil_header_pmem_init(zil_header_pmem_t *zh);
hdl_t* prb_setup_hdl(prb_t *prb, const zil_header_pmem_t *hdr);
void prb_tear_down_hdl(hdl_t *hdl,
    bool abandon_claim, zil_header_pmem_t *upd);

```

HDL recover their DRAM state from the ZIL header (`zil_header_pmem_t`). The `setup prb_setup_hdl` function takes a constant pointer to the last-synced header. The pointer is not internalized in HDL — the pointee’s lifetime may be as short as the function call. The PRB consumer instantiates all dataset’s HDLs during pool import or whenever a new dataset is created. For new datasets, the initial value for the ZIL header (state *nozil*) is set by `zil_header_pmem_init`.

When the head dataset is destroyed or the pool is exported, the PRB consumer calls `prb_tear_down_hdl` to destroy the HDL. The PRB must be freed before all of its HDL’s have been torn down.

Claiming

```

check_replayable_result_t prb_claim(
    hdl_t *hdl, uint64_t pool_first_txg,
    zil_header_pmem_t *upd);

```

After the PRB is constructed and HDLs are set up, we must *claim* log entries of all dataset’s HDLs. The corresponding function `prb_claim` is invoked by the pool import procedure for each HDL. The `pool_first_txg` is the pool’s first new transaction group. For HDLs in state *logging*, `pool_first_txg - 1` becomes the *precrash_txg*.

`prb_claim` returns an update to the ZIL header through the `upd` out-parameter. We employ this pattern throughout the entire PRB API. It is always the API consumer’s responsibility to ensure that the update is correctly persisted in the correct transaction group.

Claiming can fail, e.g., if the procedure detects a missing entry for the HDL (claiming performs a dry-run of replay internally). The API consumer defines the error handling policy. It can either abort pool import or choose to abandon the log. To abandon the log, the API consumer must first tear down the HDL using `prb_tear_down_hdl`, `abandon_claim=true`, ...) to release claims made during `prb_claim`. Then, the API consumer resets the header using `zil_header_pmem_init` and re-instantiates the HDL using `prb_setup_hdl`.

need third option to retry with an override flag that discards the unplayable part of the log

After all HDLs have been claimed the pool import procedure starts *txg sync* which writes out the header updates made by claiming in the `pool_first_txg`.

From that point on there is no need to coordinate HDL operations between different datasets.

5.14.2 Replay

```
typedef struct { ... } replay_result_t;

replay_result_t
prb_replay(hdl_t *hdl, replay_cb_t cb, void *cb_arg);

typedef struct replay_node replay_node_t;
typedef int (*replay_cb_t)(void *rarg,
    const replay_node_t *rn,
    const zil_header_pmem_t *upd);

typedef enum {
    READ_REPLAY_NODE_OK,
    READ_REPLAY_NODE_MCE,
    READ_REPLAY_NODE_ERR_CHECKSUM,
    READ_REPLAY_NODE_ERR_BODY_SIZE_TOO_SMALL,
} read_replay_node_result_t;

read_replay_node_result_t
prb_replay_read_replay_node(
    const replay_node_t *rn,
    uint8_t *body_out, size_t body_out_size,
    size_t *body_required_size);

void prb_replay_done(
    hdl_t *hdl, zil_header_pmem_t *upd);
```

When a dataset is mounted the mounting procedure must always call `prb_replay`. If the HDL is in state *nozil*, the call is a no-op. If the HDL is in state *replaying*, the function invokes the provided replay callback for each log entry that needs to be replayed in replay order. The callback must perform the following steps for crash-consistent replay for each replayed entry *E*.

1. Start a DMU transaction *tx*.
2. Load the entry *E* into a DRAM buffer using `prb_replay_read_replay_node`.
3. Apply the change encoded in *E* to the dataset.
4. Update the ZIL header to the value of **upd*.
5. `dmu_tx_commit(tx)` the DMU transaction.

Note that the callback must use `prb_replay_read_replay_node` function because `replay_node_t` is an opaque type in the PRB API. The function forces the API

consumer to do DRAM buffering and protects against online and offline data corruption internally, as discussed in Section 5.8.

`prb_replay` can fail either due to an error returned by the callback, due to an error in the scanning phase, or due to log corruption. If the error is due to missing log entries, the struct returned by the API contains the *witness* log entry (see Section 5.6). The caller decides whether to retry replay or abandon the remaining unreplayable part of the log. End of replay must be acknowledged explicitly by calling `prb_replay_done`. Abandoning the log is done using `prb_destroy_log` (next section).

retry with an override flag that discards the unreplayable part of the log

5.14.3 Writing Entries

```
bool prb_create_log(hdl_t *hdl, zil_header_pmem_t *upd);
void prb_destroy_log(hdl_t *hdl, zil_header_pmem_t *upd);

int prb_write_entry(hdl_t *hdl,
    uint64_t txg, bool needs_new_gen,
    size_t body_len, const void *body_dram);
```

After successful replay, the HDL is always in the unwriteable state *nozil*. The log writer must use `prb_create_log` to create a new log idempotently. If the HDL is in state *nozil*, the function allocates a log GUID and transitions the HDL to state *logging*. Otherwise, the HDL must already be in state *logging* and the call is a no-op. If a new log was created, the caller must ensure that the transaction group that persists the ZIL header update has synced to disk before starting to write log entries. The caller distinguishes the cases based on the function's return value.

Log writers use the `prb_write_entry` function to write log records to the HDL. The log record is treated as an opaque blob. PRB does not provide facilities to version different version of the encoding. In addition to the body (`body_dram`, `body_len`), the log writer must provide two pieces of metadata:

fix in upstream release?

txg The transaction group $T_{i_{txg}}$ of the DMU transaction T_i whose changes C_i are encoded in the log entry (Section ??). Note: for ZIL-PMEM, this value is always the same as the log record's `lrc_txg` field.

needs_new_gen Indicates whether a new generation should be started for this log entry. It is the responsibility of the caller to serialize the start of a new generation as described in Section 5.13.2. if a thread writes an entry with `needs_new_gen` is true, that thread must be the only thread executing `zilpmem_prb_write_entry` for the given HDL. The inverse is not true: if `needs_new_gen` is false, multiple threads may write entries in par-

contrary?

allel to the same HDL. Note that writers for different HDLs need not coordinate at all.

language?

Garabge Collection

```
void prb_gc(prb_t *prb, uint64_t synced_txg);
```

Whenever *txg sync* has finished syncing a txg T to the main pool it must call `prb_gc` with `synced_txg = T`. Note that unlike writing or recovering an individual log, garbage collection is a PRB-level operation. Synchronization is handled internally.

Chapter 6

Integration into ZFS

In this chapter we describe how we integrate PRB/HDL into ZFS in three steps. First, in Section 6.1, we describe how we re-architect ZFS to support different ZIL implementations (*ZIL Kinds*) at runtime (Section 6.1). Then, in Section ??, we present our approach to make ZFS's device management layer (VDEV) aware of persistent memory devices. Finally, in Section ??, we describe how we combine PRB/HDL and PMEM SLOG VDEVs into the new *ZIL-PMEM* ZIL kind.

6.1 ZIL Kinds

Coexistence with the existing ZIL and preservation of ZFS's crash consistency guarantees are two hard requirements for ZIL-PMEM (see Section 4.1.1). Our solution to both of these problems is to re-architect ZFS to support different persistence strategies for the ZIL while sharing the code and data structures that ultimately define crash consistency semantics. In order to make the integration of ZIL-PMEM seamless to the end user (goal: simple administration), the persistence strategy is the same for all datasets in a pool. The variable that determines the pool's persistence strategy is its *ZIL kind*. The following sub-sections present how we refactor ZFS to support ZIL kinds. The existing ZIL which uses LWBs for persistence becomes the first ZIL kind called *ZIL-LWB*. Note that some listings in this section already mention ZIL-PMEM to illustrate why ZIL kinds are necessary to integrate ZIL-PMEM. However, in our implementation, all refactoring steps presented in this section are separate commits that precede the introduction of the ZIL-PMEM ZIL kind.

6.1.1 On-Disk State

As described in Section ?? the ZIL(-LWB) keeps its persistent state in the per-dataset ZIL header and the LWB chain. For ZIL kinds, we change the ZIL header to be a tagged union that uses the new `zh_kind_t` enum as a discriminant. The existing ZIL-LWB header fields are moved into the `zil_header_lwb_t` type. ZIL-PMEM's ZIL header, which we describe in Section 6, is the second member of that union. We maintain the same size as the original `zil_header_t` which simplifies the migration of the on-disk format. Figure 6.1 shows the relevant C structures before and after the changes described in this paragraph.

<pre>typedef struct zil_header { uint64_t zh_claim_txg; uint64_t zh_replay_seq; blkptr_t zh_log; uint64_t zh_claim_blk_seq; uint64_t zh_flags; uint64_t zh_claim_lr_seq; uint64_t zh_pad[3]; } zil_header_t;</pre>	<pre>typedef enum { ZIL_KIND_UNINIT, ZIL_KIND_LWB, ZIL_KIND_PMEM, ZIL_KIND_COUNT } zh_kind_t; typedef struct zil_header_lwb { /* fields of zil_header_t, * without zh_pad */ } zil_header_pmem_t; typedef struct zil_header_pmem { /* introduced later */ } zil_header_pmem_t; typedef struct zil_header { uint64_t zh_kind; union { zil_header_lwb_t zh_lwb; zil_header_pmem_t zh_pmem; } zh_data; uint64_t zh_pad[2]; } zil_header_t;</pre>
------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Figure 6.1: The ZIL header structs (in DRAM and on disk) before and after the introduction of ZIL kinds.

6.1.2 Runtime State

As described in Section 2.3.1, the ZIL(-LWB) runtime state is kept in the per-dataset object `zilog_t`. `zilog_t` tracks the in-memory representation of log records (ITXs) in the `itxg` struct until they become obsolete due to `txg` sync or need to be written to stable storage by `zil_commit`. `zil_commit` drains the ITXs into the *commit list*. It proceeds by packing the ITXs on the commit list into LWBs which it writes to disk as a chain linked by ZFS block pointers.

We observe the following properties of the code that handles ITXs and `itxgs`:

- It defines the framework for ZFS's crash consistency semantics. Whereas ZPL and ZVOL code fill the body of the log records, the organization by `itxgs` and the code that assembles the commit list constraints what can be expressed as ZIL records.
- ITXs and even the commit list are independent of the LWB chain that is ultimately written to disk. The commit list merely defines the set and order of ZIL records that need to be persisted before `zil_commit` and the parent synchronous I/O system call can return to userspace.
- The interface between the ITX- and LWB-related code is limited to the *commit list* and its contents. Whereas the ITXs are not opaque to the LWB-related code (e.g. handling of `WR_NEED_COPY`), the responsibilities are cleanly separated.

check that this was defined

Given these insights we refactor the ZIL implementation (*zil.c*) as follows:

1. Move all non-ITX functions into a separate module *zil_lwb.c* and prefix them with `zillwb_`. If the function was part of the public ZIL API, add a wrapper function with the original name to *zil.c* that forwards the call to the `zillwb_` function in *zil_lwb.c*.
2. Virtualize calls to `zillwb_` functions in *zil.c*:
 - Define a struct `zil_vtable_t` that contains function pointers with the type signature of each of the `zillwb_` functions called from *zil.c*.
 - Define `zillwb_vtable` as an instance of `zil_vtable_t` that uses `zillwb_` functions as values for the respective function pointer members.
 - Add a member `zl_vtable` to `zilog_t` that is pointer to a `zil_vtable_t`.
 - Replace all calls to `zillwb_FN()` in *zil.c* with indirect calls through the `vtable`, i.e., `zilog->zl_vtable.FN()`.
3. Make non-ITX state private to *zil_lwb.c* by turning it into a subobject.
 - Move the `zilog_t` members that are only used by the functions in *zil_lwb.c* into a separate structure called `zilog_lwb_t` that is private to *zil_lwb.c*.
 - Embed `zilog_t` as the first member in `zilog_lwb_t`.

- Add member `zlv_t_alloc_size` to `zlv_t_vtable_t` that indicates the amount of memory to be allocated when allocating a `zilog_t`.
- Add a *downcast* step to the start of each `zillwb_` function that casts the `zilog_t` pointer into a `zilog_lwb_t` pointer. The cast is safe because `zilog_t` is embedded as the first member of `zilog_lwb_t`. The majority of `zillwb_` functions operate only on the `zilog_lwb_t`-private state without accessing the embedded `zilog_t`.
- Add constructor and destructor methods to the vtable that are called after allocating the `zlv_t_alloc_sized` `zilog_t`. The `zillwb_` constructor and destructors initialize and deinitialize the private members of `zilog_lwb_t`.

The end result is best described in the terminology of object-oriented programming: **`zilog_t` is an abstract baseclass** that implements the public ZIL interface as well as ITX-related functionality and defines abstract methods for persisting log records. These abstract methods must be implemented by concrete subclasses. `zilog_lwb_t` is such a subclass that implements the LWB-based persistence strategy. For ZIL-PMEM, the `zilog_pmem_t` struct which we introduce in Section 6 implements persistence directly to PMEM. Which subclass is instantiated at runtime is determined by the `zh_kind` field in the ZIL header. Figures 6.2 and 6.3 illustrate the changes described in this section.

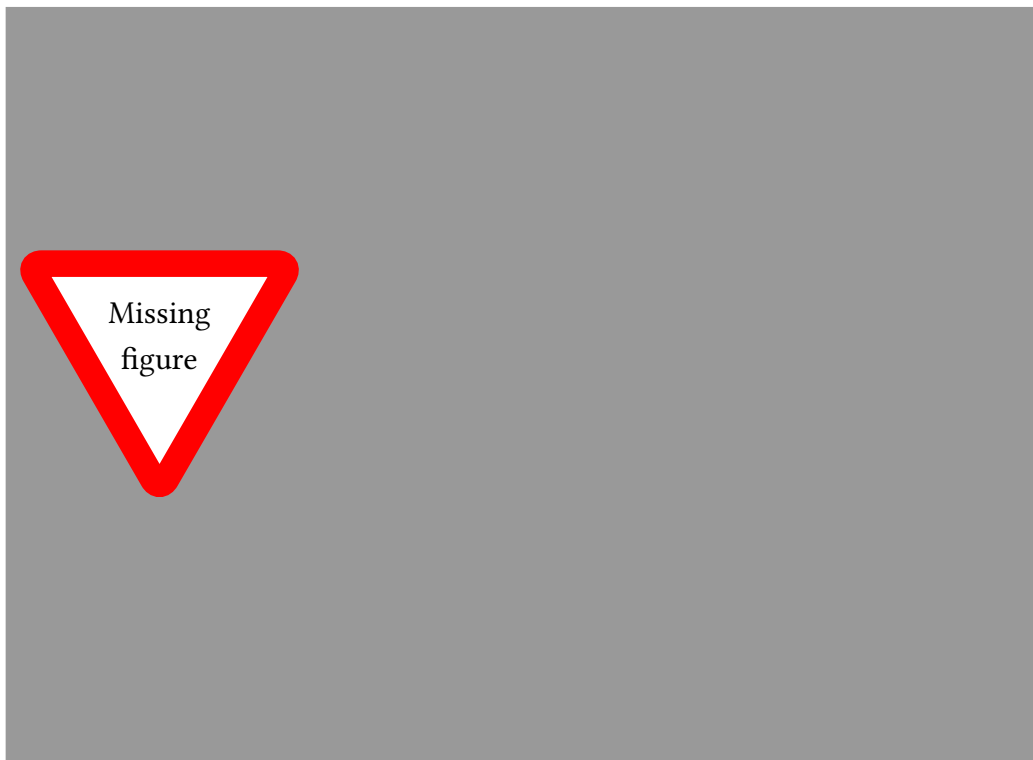


Figure 6.2: `zilog_t` before and after the introduction of ZIL kinds, visualized as a class diagram.

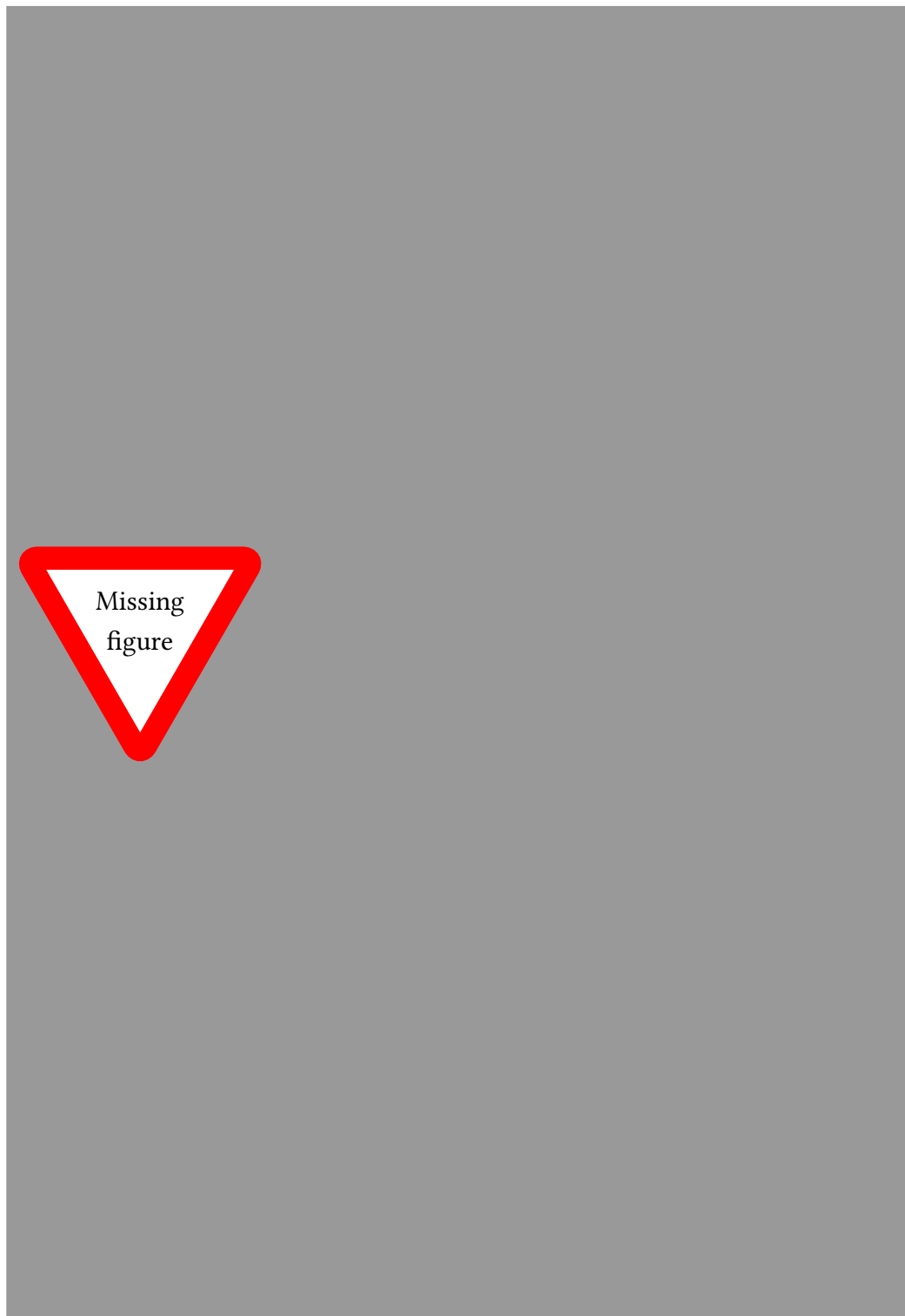


Figure 6.3: The refactoring steps described in Section 6.1.2 by example of the `zil_itx_assign()` and `zil_commit()` APIs.

6.1.3 Changing ZIL Kinds

A zpool's ZIL kind can be changed by switching over the `zh_kind` of every ZIL header in the pool. The following procedure enables online switching of the ZIL kind:

1. Ensure that all dataset's ZILs have been replayed. If not, cancel the procedure unless the caller specified to drop unreplayed logs.
2. Stop use of the ZIL API and wait until all active API calls to it have finished.
3. Wait for all ZIL entries to become obsolete by waiting for the current open txg to be synced.
4. Free all dataset's `zilog_t` instances.
5. Set all datasets' ZIL headers to the new ZIL kind's default value. The default value is the ZIL state that encodes the absence of log entries, e.g., *nozil* for ZIL-PMEM.
6. Allocate the new `zilog_t` instances bring them into a usable state. The allocation routine uses `zh_kind` to select the (new) vtable, allocates `zlvt_alloc_size` bytes for the new `zilog_KIND_t` and runs the ZIL kind specific constructor.
7. Enable access to the ZIL API by unblocking the waiting callers.

Note that it is critical that Step 5 happens atomically for all ZIL headers in the pool. Otherwise, a crash could result in a pool with mixed ZIL kinds.

Due to time constraints we have **not yet implemented** the procedure outline above. As a stop-gap solution, we added a kernel module parameter `zil_default_kind` that defines the ZIL kind that is used for the `zh_kind` field in new dataset's ZIL headers. Our implementation has no mechanism to change `zh_kind` over the lifetime of a dataset. It is possible to mix different ZIL kinds in the same pool by changing `zil_default_kind` before creating a new dataset. This works correctly but is undesirable from a usability perspective because we want ZIL kinds to be transparent to the user.

6.1.4 ZIL-LWB Suspend & Resume

The ZIL API provides the `zil_suspend` and `zil_resume` functions. `zil_suspend` pauses all ZIL activity and waits until all log entries are obsolete before returning to the caller. `zil_resume` reverts the state to normal operation. Compatibility code for versions of ZFS prior to the *fast snapshots* rely on ZIL suspend & resume for taking snapshots. For newer pool versions, the only consumer is `spa_reset_logs`: when removing a SLOG from the pool, the ZIL is temporarily suspended to ensure that the SLOG does not contain valid log entries. After the SLOG is removed, the ZIL is resumed and the metaslab allocator uses the remaining SLOGs or the main pool devices to allocated LWBs.

With regards to ZIL kinds, only the `spa_reset_logs` use case is relevant since ZIL kinds require a more recent pool version than the *fast snapshots* feature. Our design for changing ZIL kinds (Section 6.1.3) requires suspension of all ZIL activity and thus subsumes the ZIL-LWB specific `zil_suspend` and `zil_resume`. However, the compatibility code for pools before *fast snapshots* needs to be maintained in a way that does not depend on the ZIL kinds feature. Due to time constraints, we were **unable to address suspend & resume in our design**. We expect that the solution will be highly dependent on implementation-level constraints.

6.1.5 ZIL Traversal & ZDB

Whereas ZIL writing and replay are abstracted away by the `zilog_t` refactoring, there are several cases where the raw ZIL-LWB chain is traversed directly using the `zil_parse` function. `zil_parse` exposes several ZIL-LWB specific implementation details to its callers such as blockpointers and the concept of LWBs. This is problematic for ZIL kinds because not every conceivable ZIL kind uses these concepts — ZIL-PMEM being the obvious example. We investigate all users of the ZIL traversal code and come to the conclusion that there is no need for a generalized interface that every ZIL kind needs to implement. The basis for this decision is a manual audit of all `zil_parse` callers:

`dmu_traverse` This module implements a callback-based traversal of the `zpool`'s data structures. It is used to implement many ZFS features, e.g., *zfs send*. If a dataset is traversed that is a head dataset (i.e., not a snapshot) and its LWB chain has been claimed, the LWBs are included in the traversal.

`dsl_scan_zil` During a *zpool scrub* (data integrity check of the entire pool), this function traverses claimed LWB chains.

`spa_load_verify` During pool import this function uses *dmu_traverse* to validate data structures that were modified in the last synced transaction groups.

`zdb_il.c` The *ZFS debugger* interprets the ZIL header of head datasets, traverses their LWB chain, and dumps its contents to `stdout`.

Most consumers of *dmu_traverse* operate on snapshots, not head datasets, and therefore do not trigger ZIL chain traversal. The *dsl_scan* and *spa_load* code only traverses the ZIL but does not access its data — the data integrity checks that are done for validation are implemented transparently in the ZIO read pipeline that is used to load the LWBs in the ZIL chain. One compatibility code path (`old_synchronous_dataset_destroy`) uses ZIL traversal to free the ZIL blocks, but can be replaced with a more recent API (`zil_destroy_sync`). *zdb* is an exception since its whole purpose is to interpret the ZIL chain for debugging purposes.

Given this analysis, we come to the conclusion that a generic ZIL traversal API is not necessary in practice. Hence, the ZIL vtable does not include such an API. To maintain pre-ZIL-kinds behavior, we make the following changes as a precursor to the refactoring of `zilog_t` which we described in Section 6.1.2. We change the `zil_parse` API to work directly on a `zil_header_lwb_t*` instead of `zilog_t`. We also rename the function to `zillwb_parse_phys` to reflect the fact that it is specific to ZIL-LWB and does not affect runtime state. We change the `dmu_traverse` API so that callers must be explicit about ZIL-LWB traversal. To avoid ZIL-LWB specifics in the DMU traversal callback, we change `old_synchronous_dataset_destroy` to use `zil_destroy_sync` instead of DMU traversal.

It is our impression that traversal of the ZIL-LWB chain for the purpose of data integrity checking is moot. The reason is that ZIL-LWB cannot distinguish data corruption from the end of the LWB chain because it relies on invalid checksums to detect the end of the chain. It is conceivable that, for claimed-but-not-replayed ZILs, lost LWBs could be detected and surfaces as errors to the user. However, the current ZIL implementation suggests that this case has never been a top priority of the ZIL design. For example, losing a claimed-but-not-replayed log entry can leak space in the main pool or cause a crash during pool import.

Other ZIL kinds or a future revision of ZIL-LWB might be able to detect the loss of log entries and handle such situations more gracefully. In that case, an abstract integrity check method on the vtable that operates only on the physical structure, not runtime state, might be advisable.

github issues
+ what I asked
on slack today
2021-05-05

6.1.6 ZIL-LWB-Specific Callbacks

There are several callbacks in the ZIL API that are necessary for ZIL-LWB. They need to be considered for our ZIL kind refactoring.

zil_lwb_add_txg Necessary to keep the in-DRAM representation of an LWB alive when writing `WR_INDIRECT` blocks.

zil_lwb_add_block Necessary for an optimization that minimizes the amount of flush commands that are sent to the SLOG device.

zil_bp_tree_add During a ZIL traversal with `zil_parse`, this API was used to avoid doing operation more than once for a given blockpointer. It is an implementation detail of ZIL-LWB that is only a public ZIL API because it is used by `zdb`'s ZIL traversal code.

None of them are necessary for ZIL-PMEM nor are they likely to be for other ZIL kinds. Therefore, we prefix the callback functions with `zillwb_` and move them to `zil_lwb.c`. The original call sites are unreachable with ZIL-PMEM which

allows us to ignore them for the remainder of this thesis. We recommend that future work replace these statically dispatched callbacks with dynamic callbacks through function pointers to enforce decoupling.

6.1.7 Reflection & Alternatives

The general concept of ZIL kinds and the vtable-based implementation add complexity to the ZIL code. In this section we discuss the alternatives that we considered for supporting multiple ZIL implementations at runtime.

Alternative #1: We considered to move the level at which we dispatch into ZIL kind specific code one API layer upwards. This would make `zilog_t` and ITX a private implementation detail of ZIL-LWB. The API layer at which the virtual dispatch would take place are the `zfs_log_OP` and `zvol_log_OP` helper functions which create and assign ITXs for ZPL and ZVOL operations. The sequence diagram in Figure 6.4 provides an example of how they are used by a write system call. Declaring this API layer the interface for ZIL kinds would allow ZIL implementations to choose freely how they want to represent log records in DRAM and on stable storage. This additional freedom could be used by future ZIL kinds to implement a different log structure with better scalability or more fine-grained crash consistency guarantees. For example, an earlier design for ZIL-PMEM used a graph-based log structure where log entries for a file in the same dataset could be written in parallel. However, the additional freedom also has significant drawbacks that ultimately led us to the design presented in the previous subsections:

1. The `zfs_log_OP` family of functions only addresses the ZIL write path. There are no equivalent abstractions that wrap the `zilog_t` APIs for ZIL replay or traversal. In order to have a single clean abstraction at the level of `zfs_log_OP`, it would be necessary to extend this API so that `zilog_t` could be hidden as an implementation detail of ZIL-LWB. We were not confident in our ability to design this API extension without the risk of introducing a leaky abstraction.
2. ZIL kind specific functions for logging would also require ZIL kind specific replay functions. In ZIL-LWB this is a non-trivial amount of code. ZIL kinds such as ZIL-LWB that only implement a different persistence strategy would have to duplicate this code, pointlessly increasing maintenance cost.
3. We find it undesirable to allow ZIL-kinds to implement different crash consistency guarantees, in particular if ZIL kinds switch automatically depending on SLOG configuration as proposed in this chapter (Section 6.1.3). Centralizing the ITX code and forcing every ZIL kind to fit into the ITX

model is the best way to ensure that crash consistency guarantees are the same across ZIL kinds.

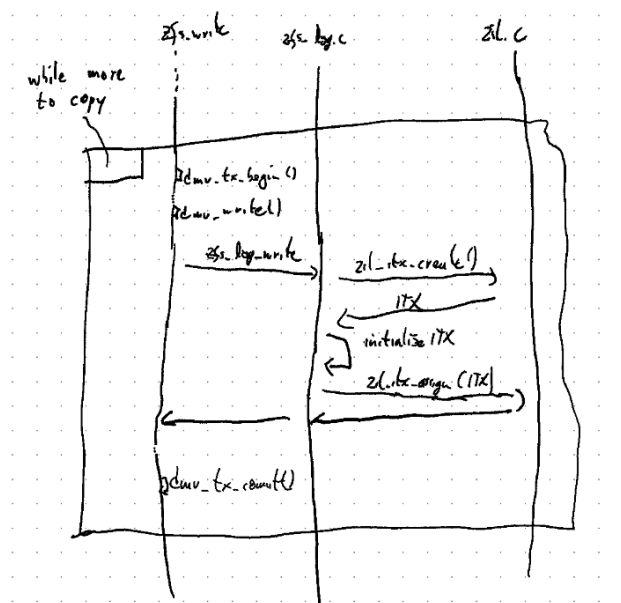


Figure 6.4: Sequence diagram of the APIs involved in creating the ITXs for a write system call.

Alternative #2: Since we identified the ZIO pipeline as the main source of latency in ZIL-LWB, we considered sharing LWBs as a concept between all ZIL kinds. In that scenario, the ZIL kinds would merely be an alternative to the ZIO pipeline. A prototype that adopts this approach has been presented at the OpenZFS 2020 Developer summit by OpenZFS, targeting NVMe drives [Ope20]. We found this layer of the ZFS software stack to be too restrictive for ZIL-PMEM:

ref chapter

1. The timeout mechanism for packing multiple entries into a single LWB would add unnecessary latency overhead. This is in conflict with one of our requirements (see Section 4.1.1).
2. This notion is shared by the database community which has deemed group commit schemes — such as LWB timeout — unfit for PMEM (see Section ??).
3. The design space for PRB/HDL would have been severely constrained. In particular, fully parallel logging to the same HDL would not have been possible because the LWB chain is inherently sequential. Our ITXG bypass for ZVOLs (see Section ??) shows how we can use this feature of PRB/HDL to increase scalability without weakening crash consistency guarantees.

ref

check,ref

We believe that our design for ZIL kinds introduces ZIL kind specific behavior at a layer that allows for shared crash consistency semantics *and* flexibility in the storage model while minimizing code duplication and thus (hopefully) maintenance cost. The interface defined by the `vtable` is a clean abstraction although some design questions (see 6.1.3 and 6.1.4) as well as some LWB-specific APIs (see 6.1.5 and 6.1.6) remain. The state of each `zilog_KIND_t` is truly private to the ZIL kind's implementation. Neither ZIL-LWB nor ZIL-PMEM access the ITX-related state in the embedded `zilog_t` directly, but only through the `zilog_t` method that computes the commit list. If requirements change in the future, the design alternatives presented in the previous section should be considered.

6.2 PMEM-aware SPA & VDEV layer

Prior to the work presented in this thesis, ZFS had no concept of persistent memory. Fortunately, the requirements of ZIL-PMEM are very limited:

- `/dev/pmem` SLOGs must be recognized as PMEM SLOG vdevs when they are added to the pool.
- The PMEM SLOG's allocatable space must be directly accessible via the PMEM programming model (load/store instructions, cache flushing, ...). If direct access is not possible, the PMEM SLOG vdev must be considered faulty and pool import must be refused. We assume that direct access cannot be lost if it is established.

future work:
memory hot-
plug

We add a new boolean attribute `is_dax` for disk VDEVs in the `zpool` config format. The attribute indicates whether the VDEV supports direct access through the Linux Kernel's DAX APIs. DAX capability is determined by the `zpool` command when creating or adding devices to a pool, using `libblkid`. When opening a vdev marked `is_dax`, the kernel module ensures that all of the block device's sectors are mappable as one contiguous range of kernel virtual address space. Failure to establish this mapping fails the onlining process, leaving the vdev in state `VDEV_STATE_CANT_OPEN`. By default, this state prevents the pool from being imported. If the VDEV is added as a log device, the import process allows the user to specify an override flag, causing the logs to be dropped. Note that the `is_dax` feature applies to all VDEVs and is independent of ZIL-PMEM. It merely records the fact that a VDEV is required to be directly accessible via the DAX APIs. This ensures that future versions of ZFS can leverage DAX capability of main pool VDEVs without needing to change the on-disk format. Consequently, `is_dax` becomes an independent *zpool* feature.

remove that
sentence? only
confusing

6.3 The ZIL-PMEM ZIL Kind

6.3.1 Activation

The ZIL-PMEM ZIL kind is used instead of ZIL-LWB based on the following rule:

If the pool has exactly one SLOG and that SLOG is `is_dax`, the ZIL kind is ZIL-PMEM. Otherwise, it is ZIL-LWB.

We evaluate this rule during pool import and whenever the pool config changes, i.e.,

- during pool creation,
- on `zpool add` and `zpool remove`, and
- when `zpool import` is told to drop log devices via the `-m` flag.

To determine whether we need to change ZIL kinds, we compare the resulting desired ZIL kind $k_{desired}$ with the `zh_kind` value in any dataset's ZIL header. If they do not match, we change the ZIL kind as described in Section 6.1.3.

The idea behind this approach is to keep administration simple which is a defining paradigm for ZFS [Bon+03]. From the system administrator's perspective, ZIL-PMEM is automatically used if it supports the user-specified SLOG device configuration, i.e., a single leaf `/dev/pmem` device. Otherwise, the system transparently falls back to ZIL-LWB.

Note that this design makes the rule that decides the pool's ZIL kind effectively part of the `zpool`'s on-disk format. For better maintainability and forwards-compatibility, the pool's ZIL kind should be encoded in the `zpool` config instead.

6.3.2 PRB Construction

If the pool's ZIL kind is determined to be ZIL-PMEM, ZIL-PMEM uses the entire allocatable space of the (single) PMEM SLOG VDEV to construct the pool's PRB instance on top of it. The pool import procedure allocates the `prb_t` and stores the pointer in the DRAM object that represents the imported pool (`spa_t`). When the pool is exported, the `spa_unload` procedure frees the `prb_t`. The `prb_t`'s lifetime is thus a contiguous sub-span of the lifetime of `spa_t`.

Immediately after allocating the PRB, pool import sets up the chunk objects and adds them to the PRB. If the pool is being created, we reset the chunks to a known PMEM state (`prb_chunk_initialize_pmem`) before adding them to PRB. Due to time constraints, our implementation uses a simplistic partitioning scheme that does not need to persist any metadata:

correct tense?

impl is a little different, doesn't matter

- Chunks have a hard-coded chunk size of 128 MiB.
- We partition the allocatable space as a contiguous array of PMEM segments of 128 MiB, starting at offset zero.
- Assuming A bytes of space, this yields $nchunks := A \gg 27$ chunks and less than 128 MiB of wasted space at the end of the allocated space.
- We do not store $nchunks$ anywhere. Instead, we prohibit online resizing of the PMEM SLOG vdev which allows us to deterministically re-compute the value.

6.3.3 Dataset & HDL Lifecycle Synchronization

language

don't mention MOS here, we'd need to explain it

we dropped the name a few times in the ZIL kinds section, ok to use here without explaining?

Whereas `prb_t` and `spa_t`'s allocation lifecycles line up nicely, the same is not true for datasets and HDLs: HDLs must be set up during pool import and live until either their dataset destroyed or the zpool is exported. In contrast, the per-dataset runtime state (`dsl_dataset_t`, its `objset_t` and the `objset`'s `zilog_t`) is allocated *on demand* when its consumer *holds* it by its ID. If the consumer is the first holder, the *hold* procedure performs the allocation and recovers state from the corresponding persistent structure (`dsl_dataset_phys_t`, `objset_phys_t` and, partially, `zil_header_t`). It then increments a reference counter and returns the pointer to the runtime object to the consumer. Conversely, if there was already another holder, the allocation and recovery steps are skipped. Once a consumer is finished with a dataset, it *releases* its hold on the DRAM structure. The last consumer that releases the structure frees the DRAM allocation.

The mismatch of HDL and dataset object lifetimes is relevant to ZIL-PMEM because `zilog_pmem_t` needs access the corresponding HDL for claiming, replay, and writing log entries. Our solution is to add a search tree to `spa_t` which we call *HDL map*. It maps from the dataset's object set ID to its HDL. We add the necessary callbacks to synchronize the *HDL map* with the dataset layer.

- During pool import, after setting up the PRB but before claiming, we iterate over the head datasets in the pool, set up the HDLs for each of them, and add them to the *HDL map*.
- When a dataset is created (`dmu_objset_create_impl_dnstats`), we set up its HDL and add it to *HDL map*.
- When a dataset is destroyed (`zil_destroy_sync`), we tear down the HDL and remove it from the *HDL map*.

Whenever the ZIL-PMEM implementation needs access to the dataset's HDL, it looks up the HDL in the *HDL map* and acquires a reference to it. We use reference counting to assert that there are no dangling references when the HDL is destroyed in `zil_destroy_sync`. We protect the *HDL map* against concurrent

modifications using ZFS's *read mostly lock* which is a reader-writer-lock that is optimized for read locks. Figure 6.5 visualizes the constellation of objects, their lifetimes, and reference relationships. language



Figure 6.5: The different allocation lifecycles of HDLs and datasets, and how we bridge them for ZIL-PMEM.

6.3.4 ZILOG_PMEM_T

The `zilog_pmem_t` structure and its methods in the `zilpmem_vtable` implement the ZIL-PMEM ZIL kind. Its role is that of an adaptor between the shared `itxg` structure that is shared among all ZIL kinds and the HDL that is associated with the dataset via the *HDL map*. Apart from the exceptions listed below, the methods in the vtable are thin wrappers around HDL's interfaces. The general pattern for a ZIL-PMEM method is as follows:

1. Acquire a reference to the dataset's HDL from *HDL map*.
2. Invoke the HDL method.
3. Release the HDL reference.

To avoid reference counting overhead for HDLs on the write path, we acquire a HDL reference once when the ZIL is opened for writing and hold it in the `zilog_pmem_t` until the dataset is unmounted.

The following cases required additional code to adapt between the two domains:

Error Handling During Pool Import & Claiming The ZIL vtable was extracted from the original ZIL(-LWB) and inherits its model for integrity checking

and claiming of the ZIL during pool import. In that phase, ZIL-LWB checks the ZIL chain twice: the first pass (`zil_check_log_chain`) checks for log corruption to fail the pool import early if the log is corrupt. The second pass (`zil_claim`) does ZIL-LWB's version of claiming and records the maximum claimed log record in the ZIL header. It discards traversal errors because it assumes that no online log corruption happened since the first pass. (Note that `zil_claim` does have a return value, but only due to software-technical reasons. It must always return zero.) ZIL-LWB's approach opens a window for time-of-check vs time-of-use bugs which we discovered and reported during development. In contrast, PRB/HDL does checking and claiming in a single pass (`prb_claim`), reports corrupted log state through its return value, and is designed to correctly handle online log corruption during claiming and replay (see Section 5.8).

ref

Due to time constraints, we have not refactored the zpool import procedure to allow claiming to fail. Until that shortcoming of our implementation is addressed, we trigger a kernel panic if `prb_claim` returns an error.

ZIL Header Updates Remember from Section 5.14.1 on the PRB/HDL API that the responsibility of *persisting* the ZIL-PMEM header is split: HDL APIs that update the ZIL header return the updated version through an *out*-parameter and the API consumer is responsible for persisting the update to the main pool.

Remember also the `zil_sync` API from Section 2.3.1. *Txg sync* invokes this method on every dataset's `zilog_t` in the *syncing* txg. The role of `zil_sync` is to modify the dirty copy of the ZIL header that is stored in `objset_phys_t::os_zil_header`. Unlike most components at this software layer of ZFS, the ZIL implementation itself is responsible for queuing up changes that need to be applied for the *open*, *quiescing* and *syncing* txgs.

ZIL-PMEM queues the updated header values returned by the HDL APIs in a 4-ary array that is indexed by `[txg&3]`. Each cell contains the tuple (`txg: u64, header: zil_header_pmem_t`). A header update (`txg_u, header_u`) overwrites the cell at index `txg_u&3`. Updates must only be made from the *open* or *quiescing* txgs. `zil_sync` then reads cell `C := updates[syncing_txg&3]`. If `C.txg == syncing_txg`, it updates the ZIL header to `C.header`. Otherwise, it does not modify the ZIL header.

We implement `zil_commit` as follows:

1. ~~Acquire~~ `zil_commit` acquire a `zilog_pmem_t`-wide mutex.
2. Get the *commit list* from the *itxg* data structure.

3. Invoke the HDL's `prb_write_entry` method for each log entry in the commit list, starting a new generation for every entry.
4. Release the mutex.

The mutex ensures that committers are serialized. This is critical for the correctness of *sync* or *fsync* which are cumulative. Assume two threads *A*, *B* that issue the *sync* system call. Assume *A* starts a *sync* system call and drains all ITXs from the *itxg* into its commit list. If *A* is now preempted by *B* which also starts a *sync* system call, *B*'s commit list will be empty. *B* would return to userspace, giving the caller the impression that all dirty data has been synced, although they still need to be written to the HDL by *A*.

Starting a new generation for every entry ensures that the commit list's entries will be replayed in commit list order. Note that the `zilog_pmem_t`-wide mutex already serializes the start of the new generations as required by the HDL API, see Section 5.14.3.

Serializing the entire `zil_commit` call allows for less parallelism than ZIL-LWB's *commit ITXs*. Commit ITXs implement a sort of pipelining by allowing log writers to issue LWB ZIOs in parallel while ensuring that they only return from `zil_commit` after all LWBs up to and including the log writer's last LWB have been written. Whereas a similar approach could be applied to ZIL-PMEM in principle by mapping parallel writers to the same generation, the current implementation of *commit ITXs* is too tightly coupled to the concept of LWBs and the ZIO pipeline.

language, help

WR_NEED_COPY Chunking Remember from Section 2.3.1 that the ZIL allocates ITXs for *all* changes to a dataset in DRAM. The ITXs are queued in the *itxg* structure from where they are either `zil_committed` or freed when the *txg* syncs. To avoid unnecessary memory usage and performance overhead, ITXs that log *write* operations do not always contain a copy of the written data. Instead, the `zfs_log_write` function creates a `WR_NEED_COPY` ITX. Such ITXs only contain the object number of the DMU object and the range that was modified. The creation of the log record is deferred until `zil_commit` which reads the range from the DMU object and split it up into many *write* log records.

ref it, we introduce it in the alternative approach for ZIL kinds section

For ZIL-PMEM, we implement an isolated structure that turns a single `WR_NEED_COPY` record into an iterator over *write* log entries. `zil_commit` writes each entry yielded by the iterator to the HDL using `prb_write_entry`. Whereas the repeated chunk acquisition and dependency tracking during each HDL

write incurs a small overhead, it also increases fairness among HDLs if the commit slots are contended.

6.3.5 ITXG Bypass For ZVOL

ZIL kinds are forced into the *commit list* model defined by the shared *itxg* structure which leads to an inherently sequential representation of the ZIL contents: ZIL-LWB is a long chain of log entries grouped in LWBs, and ZIL-PMEM starts a new generation for every entry on the commit list to encode it properly. We believe that this architecture is the best compromise for consistent behavior, performance, code duplication, and maintainability (see Section 6.1.7). Our evaluation in the next chapter shows that ZIL-PMEM yields significant latency benefits and can saturate PMEM bandwidth when performing synchronous I/O from multiple threads to *separate* datasets. However, the benchmarks show that multi-thread scalability within a *single* dataset is limited. To demonstrate that PRB/HDL is not the bottleneck, we develop an *itxg bypass* mode for ZIL-PMEM which we present in this section.

OpenZFS Background: ZVOLs

The mode only works for ZVOLs, which are sparsely allocated virtual block devices backed by the zpool. Other block device consumers, e.g., virtual machines or other file systems, can treat the ZVOL as any other block device exposed by the kernel. ZVOLs are implemented as a dataset with a single DMU object that contains the virtual block device's data. When a block device driver accesses the ZVOL block device (*read*, *write*, *discard*), ZFS maps the block device operations to DMU operations on the object. The modifications are logged to the ZIL as ITXs. If the block device operation has synchronous semantics, ZFS calls `zil_commit` before acknowledging completion of the block device operation. The following pseudo-code describes how the ZVOL processes a block device operation.

```

Input:
  op      block device operation
  zv      zvol
Steps:
  if op.pre_flush:
    zil_commit(zv.dmu_object_id)

  match op {
    Read(...) => {
      dmuf_read(zv.dmu_object, op.buf, ...)
    },
    Write(...) => {
      tx := dmuf_transaction()
      dmuf_write(tx, zv.dmu_object_id, ...)
      itx := zil_itx_create(...)
    }
  }

```

```

        zil_itx_assign(itx)
        dmu_tx_commit(tx)
    },
    Discard(...) => {
        // analogous
    }
}

if op.post_flush:
    zil_commit(zv.dmu_object_id)

op.done() // acknowledge completion

```

The Linux kernel's model for block device I/O is asynchronous. The assumption is that the block device IO is submitted (`submit_bio`), processed asynchronously by the storage device, and eventually marked as completed (`BIO_END_BIO`). By default, ZVOL uses a thread pool (*taskqs*) to process block device I/O asynchronously and in parallel. The `zvol_request_sync=1` tunable changes this behavior to synchronous mode where the procedure outlined above is executed in synchronously in `submit_bio`. We will refer to this tunable in the Section ?? where we evaluate the ITXG bypass.

ITXG Bypass

In *itxg bypass* mode, we do not queue ITXs in *itxg* but write the log entry/-ies directly to the HDL instead. The role of `zil_commit` is reduced to starting a new generation which is necessary to encode a conservative approximation of the logical dependencies between block device operations. Remember from Section 5.13.2 that the thread that writes the first entry of a new generation must have exclusive access to the HDL. We use a reader-writer-lock to achieve this behavior, as illustrated by the following pseudo-code. Note that the actual implementation of the bypass mode is contained within `zil_commit` and `zil_itx_assign`.

language...?

```

Input:
  op      block device operation
  zv      zvol
  rwl     reader-writer-lock
  start_gen

Steps:
  if op.pre_flush || op.post_flush:
    zv.rwl.write_lock()
    zv.start_gen = true
  else:
    zv.rwl.read_lock()
    if zv.start_gen:
      zv.rwl.upgrade()

  match op {
    Read(...) => { ... }
    Write(...) => {
      tx := dmu_transaction()
      dmu_write(tx, zv.dmu_object_id, ...)
    }
  }

```

```
        itx := zil_itx_create(...)
        prb_write_entry(itx.into_log_entry(), needs_new_gen=zv.start_gen)
        assert zv.start_gen => zv.rwl.holding_write_lock
        if zv.start_gen:
            zv.start_gen = false
            dmu_tx_commit(tx)
    },
    Discard(...) => { /* analogous */ }
}

if op.post_flush:
    assert zv.rwl.holding_write_lock
    zv.start_gen = true

op.done() // acknowledge completion
```

Chapter 7

Evaluation

In this chapter we evaluate whether ZIL-PMEM meets the project goals established in Section 4.1.1.

7.1 Usability & Architecture

Our high-level requirements for ZIL-PMEM were: simple administration, sharing of PMEM as a pool-wide resource, same crash consistency guarantees, and coexistence with the existing ZIL(-LWB). Our design meets all of these requirements:

Simple Administration Activation of ZIL-PMEM is completely transparent to the administrator. There is no change in user experience. (See Section 6.3.1 for details.)

Pooled Storage PRB wraps the PMEM SLOG device and shares it as a pool-wide resource among all HDLs / datasets. (See Section 6.3.2 for details.)

Coexistence The introduction of *ZIL Kinds* allows for coexistence of ZIL-LWB and ZIL-PMEM in code and at runtime. The layer at which ZIL kinds were introduced in the architecture allows for sharing of all code that deals with the ZIL's logical structure while providing sufficient freedom to implement a more efficient persistence strategy.

Same Guarantees `zilog_pmem_t` maintains the same crash consistency guarantees as ZIL-PMEM, courtesy of the shared logical structure and PRB/HDL's guarantees.

7.2 Correctness

We use unit tests and integration tests to validate that our implementation handles expected scenarios correctly.

7.2.1 PRB

We test PRB/HDL’s functionality in user space which is possible because PRB/HDL is included in the *libzpool* user-space library. *libzpool* contains the majority of the ZFS kernel module’s code — it is used to implement *zdb* (the ZFS debugger) and the *ztest* stress testing tool. We implement our unit tests in Rust to leverage its expressive type and macro system, rich standard library, and its built-in testing harness. We use the popular *bindgen* crate to generate the bindings to *libzpool* and implement a few idiomatic wrappers to reduce boilerplate code.

API Walkthrough

We codify the walkthrough of the PRB/HDL API that we presented in Section 5.14 in a large test:

1. Allocate a ZIL header on the stack.
2. Create a chunk whose space is allocated from the heap.
3. Construct PRB and add the chunk to it.
4. Setup a HDL.
5. Create a log for the HDL.
6. Write two entries for $txg = 2$ body values 23 and 42 to the HDL.
7. Teardown the object set.
8. Destroy the PRB with `free_chunks=false`. This moves ownership of the chunk moves back from PRB to us.
9. Construct a new PRB instance with the same chunk.
10. Setup a HDL from the ZIL header.
11. Trigger claiming.
12. Trigger replay, and record the replay callback invocations. For each invocation, we read the body length and content, assert that no read error occurred. We record body content and the ZIL header update as records in a list. We report successful replay for every callback invocation.

13. After replay is complete, we compare the recorded list's contents to the expected replay order.

7.2.2 ZIL-PMEM

Correct Handling Of Obsolete Entries

PRB/HDL assumes that there are only three unsynced txgs at any given time. This manifests in frequent use of the `txg&3` indexing idiom that is common in ZFS. Whenever state needs to be kept for each of the unsynced txgs, a 4-ary array is used to represent it. The `txg&3`'th element in the array then contains `txg`'s state. The `txg` is most often repeated within the per-`txg` state to detect when an array element is re-used.

In the context of PRB/HDL, it is critical for correctness that an entry E_i for `txg` $E_i.txg$ is only written if it is either newer than txg_{open} or one of $\{txg_{open}, txg_{open} - 1, txg_{open} - 2\}$ (dependency tracking, see Section 5.6). We have tests that the `prb_write_entry` API does not allow writing entries that do not meet this criterion. We also cover the replay path and ensure that, if the implementation were incorrect, the replay code would identify the problem with a distinct error code.

Replay Algorithm

We structure the PRB/HDL implementation such that it becomes possible to test the core replay logic that we described in Sections 5.5, 5.6, 5.7 and 5.8: The idea is to isolate the core logic in a function `replay_resume()` that takes as input

NodeSet the set of entries, represented as *replay nodes*, that were discovered for the HDL,

ReplayState a pointer to the DRAM representation of replay state,

Callback the callback that receives the *replay node* to be replayed and a pointer to the replay state that needs to be persisted to the ZIL header during replay.

It then constructs the replay sequence from *NodeSet* and invokes the *Callback* for each entry that needs to be replayed, with the replay node and ZIL header update as arguments. `prb_claim()` uses `replay_resume()` for ZIL headers in state *logging* to determine E_{seal} , and for a dry-run of replay for headers in state *replaying* to detect missing entries early (see Section 5.8, Step 1). It provides its own *Callback* (`prb_claim_cb()`) which only checks that all replay nodes can be read. `prb_replay()` uses `replay_resume()` to perform actual for the actual replay. Its *Callback* (`prb_replay_cb()`) serializes *ReplayState* to the representation stored

in the ZIL header, builds the ZIL header update, and invokes the user-provided callback. Figure 7.1 visualizes this architecture.

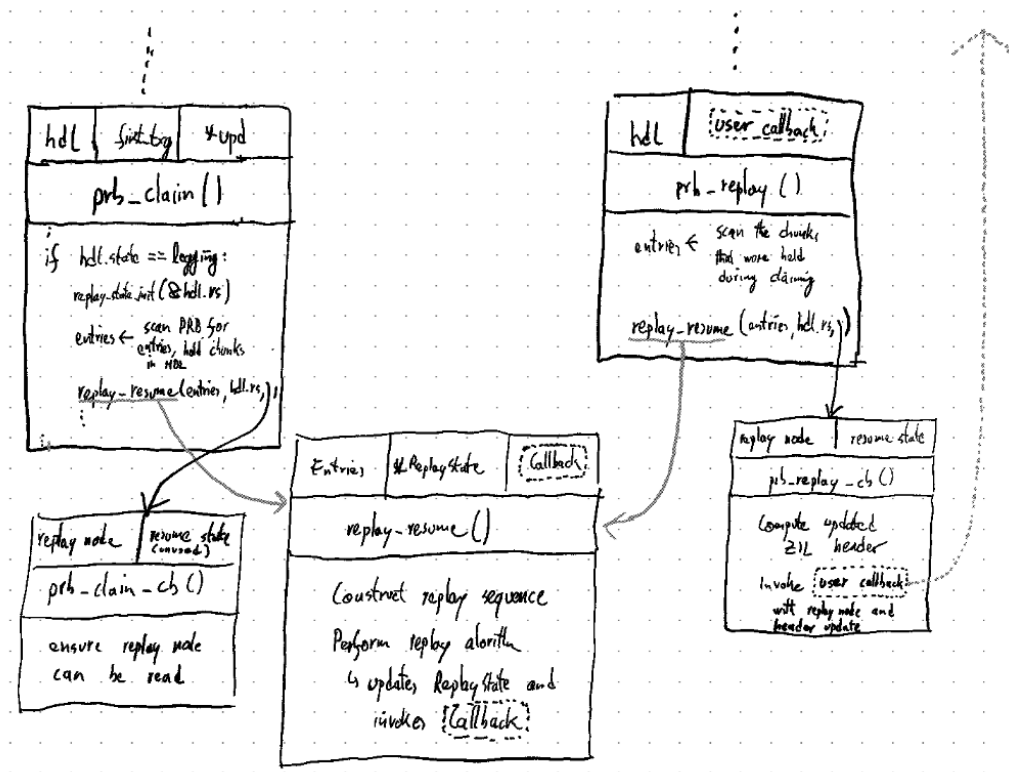


Figure 7.1: The architecture that enables code re-use and testability for the core replay logic.

`replay_resume()` pushes the responsibility to discover entries and to persist replay state to the caller. This provides maximum flexibility for testing the core logic:

Independence of PRB & Chunks Testers do not need to mock chunks or parts of PRB that would have to be scanned. They also do not need to allocate the actual entries. It is sufficient to allocate arbitrary *replay nodes* and to put them into the *EntrySet*

Independence of ZIL-PMEM It is not necessary to mock or substitute for any behavior in ZIL-PMEM.

Crash Consistency Testing To test crash consistency and correctness for resumed replay, the tester can forge an arbitrary replay state in DRAM and invoke `replay_resume()` with it.

We leverage Rust’s macro system to define test cases in a concise and expressive manner. Our goal was to minimize the mental step from the two-dimensional grid visualization (See 5.3) to a test case. The following code snippet is adapted from the test suite.

```
TestSet {
  title: "I-shape",
  entries: maplit::btreemap! {
    // SynthEntry(txg,gen,gsid, (<counters>))
    "A" => SynthEntry(3,10,1, ([ (0,0), (0,0), (0,0) ])),
    "B" => SynthEntry(4,10,2, ([ (0,0), (0,0), (0,0) ])),
    "C" => SynthEntry(5,10,3, ([ (0,0), (0,0), (0,0) ])),
    "D" => SynthEntry(4,11,1, ([ (5,1), (4,1), (3,1) ]))
  },
  tests: vec![
    test! {
      "tail_truncation_ok",
      claim_txg = 1,
      // hide entry D during replay and assert that
      // replay does not complain
      stages = stages!(single, hide=&["D"], check=OK,),
      // expected replay callback invocations
      expect_replay = vec!["A", "B", "C"],
    },
    test! {
      "'A_missing'_detected,_but_remainder_of_its_gen_is_replayed",
      claim_txg = 1,
      // hide entry A during replay and ensure that replay
      // complains about missing entries
      stages = stages!(single, hide=&["A"], check=M_ENTRIES,),
      // expected replay callback invocations
      expect_replay = vec!["B", "C"],
    },
  ],
}
```

We test the following properties:

Shapes We test replay for different “shapes” in the 2-D grid, e.g., a “V”-shape with entries for different txgs, or the “I-” shape from the snippet above.

Claim Txg We ensure that log entries from the *precrash_txg* = *claim_txg* − 1 are ignored during replay.

Missing Entry Handling We test the behavior for missing entries in the already replayed and still to be replayed parts of the replay sequence, as well as in the last generation.

Resumability We test resumability of `replay_resume()` for the cases for the cases of missing entries, re-appearing entries, and no change of *EntrySet* between resume attempts.

Entry Reappearance We test that re-appearing entries ordered before the last-replayed entry are ignored, and re-appearing entries in the unreplayed part are discovered and replayed.

Crash Tests

We make heavy use of runtime assertions in the implementation to protect against implementation errors. On the recovery path, almost all assertions are enabled for release builds of the module. On the write path, most are only enabled for debug builds to improve code generation and runtime performance. However, *libzpool* is always compiled in debug mode. We modify *libzpool* such that our test suite can install a hook whenever a runtime assertion triggers. This allows us to implement *crash tests*. At this time, we use them primarily to ensure that the PRB API triggers assertions if used incorrectly:

- Operations on a HDL in the wrong state, e.g., replay of a HDL which was not claimed.
- Attempts to set up a HDL that has already been set up.
- Writes that exceed the maximum allowed chunk size.
- Correct interplay of garbage collection and PRB during PRB allocation (critical for pool export).

7.2.3 ZIL-PMEM

Our changes to ZFS are covered (sparsely) by the testing infrastructure of the OpenZFS project.

ZFS Test Suite (ZTS) The ZFS Test Suite is an suite of automated integration tests that asserts behavior of the kernel module and CLI tools. ZTS is implemented in Korn shell and consists mostly of functional test.

... of mostly
functional
tests?

ztest A user-space binary that links against *libzpool* and stress-tests the core components of ZFS that are shared among all support platforms. The strategy is to imitate a zpool with lots of concurrent activity initiated from many threads. The focus is on code paths in the administrative layer, e.g., dataset creation & deletion or VDEV management. Whereas these code paths are also executed by ZTS, the lack of concurrency in ZTS hides potential race

conditions that are discovered using *ztest*. *ztest* is designed for easy extensibility.

ZIL Kinds

ztest has several functions that induce ZIL data path activity (*ztest_freeze()*). Whereas these functions use the public ZIL API and required little to no modification, the *ztest* mocks for datasets (*ztest_ds_t*) and its assertions of ZIL state are specific to ZIL-LWB. Due to time constraints, we

- disabled *ztest_freeze()* for ZIL-PMEM and
- changed the ZIL-LWB specific assertions to perform a checked downcast of *zilog_t* to *zilog_pmem_t*.

. In the future, this change will need to be reverted and a maintainable design for testing multiple ZIL kinds be devised. We let *ztest* run once for one hour and did not experience a failure.

ZIL-PMEM

After the test run with the ZIL kinds patch, we hard-coded */dev/pmem0* as a SLOG device in *ztest* to get basic testing exposure for ZIL-PMEM. Despite the disabled *ztest_freeze()*, this change helped us discover and fix bugs in the lifecycle of ZIL-PMEM, i.e., dataset creation, mounting, unmounting, and destruction.

Regarding the ZFS test suite, we have identified the following functional tests as relevant for the ZIL and/or SLOG devices:

proceed with
next section

TODO TODO

The following changes need to be covered by additional ZTS tests:

- As described in Section 6.3.1, we plan to automatically activate ZIL-PMEM based on the whether the SLOG device is PMEM. However, we have not yet implemented the support for online switching of ZIL kinds for a dataset (Section 6.1.3). That work should be accompanied by corresponding additions to *ztest*.
- Although we force the replay callback to use the *prb_replay_read_replay_node* function which handles *machine check exceptions (MCE)* through *memcpy_mcsafe()*, we should somehow ensure that such errors are actually handled.

7.3 Performance

We evaluate the performance of ZIL-PMEM using an extensive set of benchmarks. We seek to answer the following questions:

- What is the speedup that ZIL-PMEM achieves over ZIL-LWB in the 4k synchronous random write workload?
- What benefit does ZIL-PMEM provide to applications that frequently perform synchronous I/O?
- How does ZFS with ZIL-PMEM perform compared to other Linux filesystems that were adapted to PMEM?
- Is ZIL-PMEM a viable alternative to Linux filesystems deployed on top of with dm-writecache?
- To what degree do ZVOLS benefit from ZIL-PMEM?
- What is the benefit of the ITXG bypass for ZVOLS?
- Is PRB's *commit slots* mechanism effective at improving CPU efficiency if the maximum PMEM write bandwidth is exceeded? And what is its overhead?

7.3.1 Setup & Methodology

7.3.2 4k Synchronous Random Write Workload

The first pillar of our performance evaluation is the *fio* 4k synchronous random write workload that motivated the design of ZIL-PMEM. We already described the workload characteristics in detail in Section 3.1. We use the same zpool layout as for ZIL-LWB (30 striped NVMe partitions & 1 non-interleaved Optane DIMM region / *fsdax* namespace as SLOG). We configure PRB with *ncommitters*=3, which is the most CPU-efficient settings for a single Optane DIMM as we will show in Section 7.3.8.

Figure 7.2 shows the cumulative IOPS and per-IOP average latency measured by *fio* for 1–8 threads. With a single thread, ZIL-PMEM achieves 128k IOPS which is a speedup of 8 over ZIL-LWB (16k IOPS). ZIL-PMEM scales almost linearly to 400k IOPS at four threads where the speedup over ZIL-LWB is still 5.5. Performance does not increase further for higher thread counts. The constant offset to the *fsdax* curve for thread count 4–7 suggests that performance is limited by PMEM write bandwidth. (Remember that *fsdax* shows the raw */dev/pmem* block device's performance under the same workload. The slight decline to 346k IOPS at eight threads correlates with a similar decline in the *async* curve but not the *fsdax* curve, suggesting a CPU bottleneck (8 cores per socket) or scalability bottleneck in ZFS.

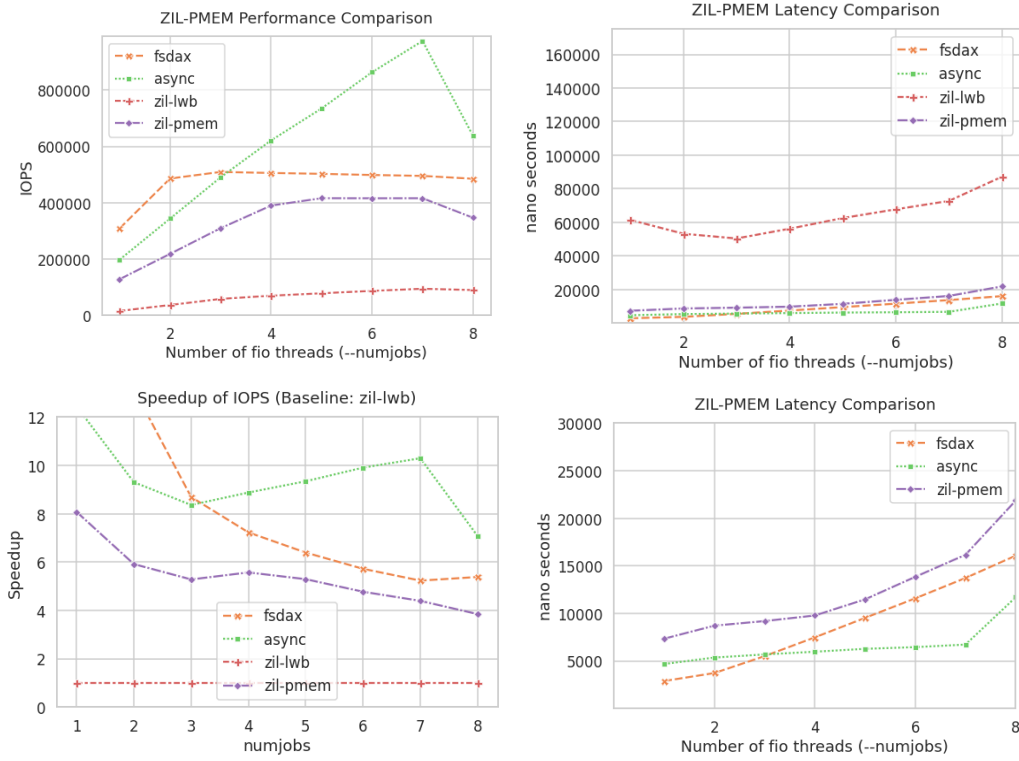


Figure 7.2: Mean IOPS and latency measured by *fio* for our 4k synchronous random write workload, by number of threads.

We observed that the speedup for `numjobs=1` and `2` varies significantly between benchmark runs but remains stable for higher values of `numjobs`. We investigate this issue by computing the coefficient of variation (CoV) of the different configurations based on the mean and standard deviation of IOPS reported by *fio*. The results are displayed in Figure ???. The *zil-pmem* CoV remains close to the CoV of *async* until `numjobs=4` from where it starts to decline towards the CoV of *fsdax*. This phenomenon correlates with the stop of increase in IOPS at `numjobs=4`, supporting our assumption that *zil-pmem*'s behavior is dominated by the PMEM hardware for `numjobs 4–7`. In contrast, *zil-lwb*'s CoV is 5x that of *zil-pmem* for `numjobs=1` and 2x for `numjobs=2` before it starts to align closely with *async*. In absolute terms, for `numjobs=1`, the standard deviation for ZIL-LWB is 4700 with a mean value of merely 16k IOPS. In contrast, for ZIL-PMEM, the standard deviation is 7557 at 128k IOPS.

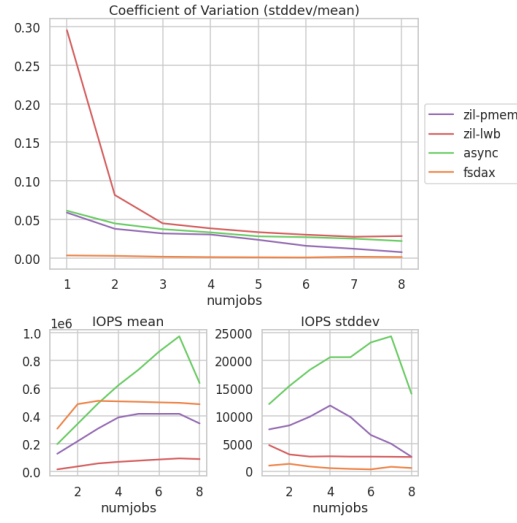


Figure 7.3: Comparison of the coefficient of variation ($\frac{stddev}{mean}$). The smaller plots display the absolute values.

To get an idea of ZIL-PMEM’s impact on tail latencies, we compare the 5th, 95th, 99.9th, 99.9th, 99.99th and 99.99th’s percentiles for completion latency, as reported by fio. *zil-pmem* follows *async* with a near-constant offset for all percentiles until the expectable knee at numjobs=4. *zil-lwb*’s 99.9th and 99.99th percentile show steadier but steeper growth with rising numjobs.

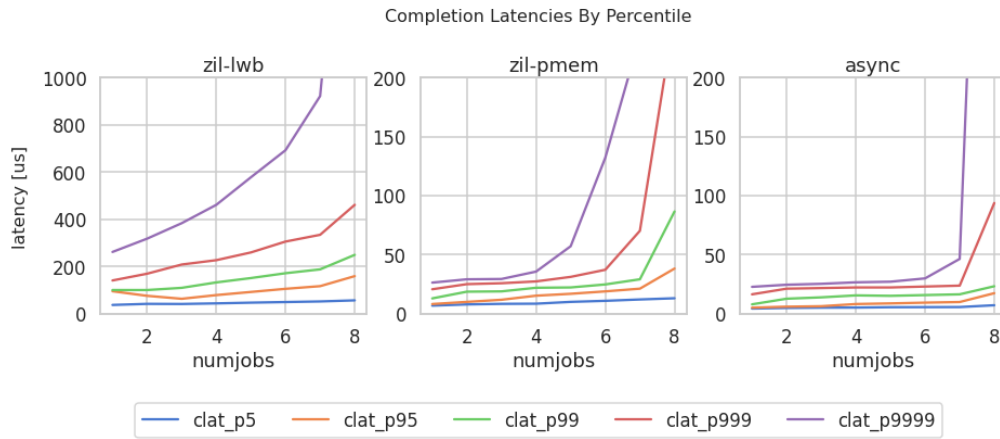


Figure 7.4: Comparison of completion latencies by percentile. Note the different y-axis scales.

We generalize our eBPF instrumentation from Section 3.3 to support ZIL-PMEM and compare the resulting breakdown of per-IOP latency in Figure 7.5. The new `zil_persistence` component is defined as the overall time spent in `zil_commit` minus the time spent filling the commit list. Whereas it is the dominant factor for latency in ZIL-LWB (80% of average per-IOP latency, as observed in Section 3.3), ZIL-PMEM only spends approximately 25% on persistence to PMEM. For ZIL-PMEM, the most dominant component is the `async` code, i.e., modification of DMU objects. And in contrast to ZIL-LWB, the shared ITX code is a noticeable component at 10% contribution to overall latency.

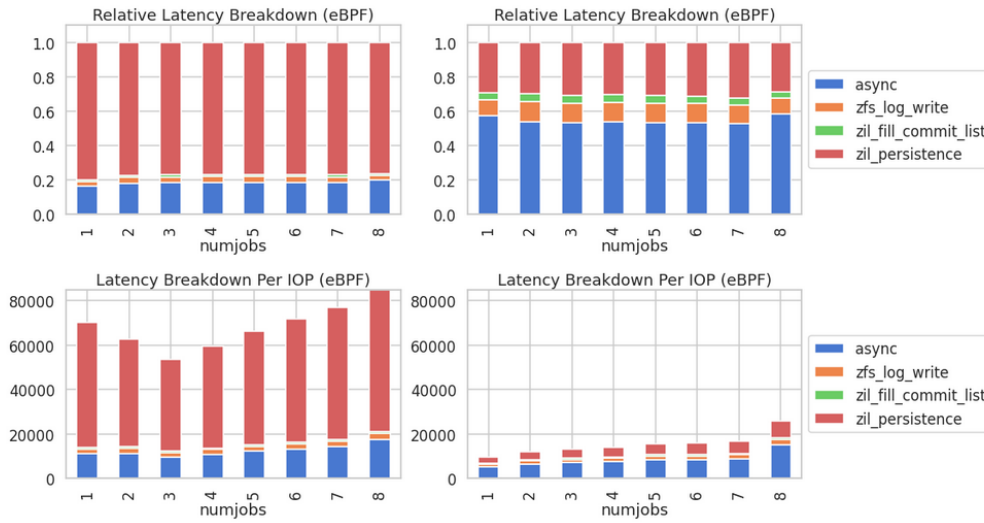


Figure 7.5: Comparison of relative and absolute latency contribution of functions executed by fio threads when they perform synchronous writes.

7.3.3 Application Benchmarks

We determine ZIL-PMEM’s impact on real-world application performance through additional macro-level and application benchmarks. We execute each benchmarks on top of different *storage stacks* that represent alternatives to ZIL-PMEM and compare the results in the subsequent sections.

We define a *storage stack* as a configuration of storage hardware and software that exposes a filesystem at a mountpoint M . Executing a benchmark *on top of* a storage stack means that the benchmark is configured to place all of its data within the filesystem mounted at M . We define the following storage stacks.

zfs-{lwb,pmem,sync_disabled} A single ZFS dataset created on a zpool with the familiar hardware configuration (30 striped NVMe partitions

and `/dev/pmem` as SLOG). For *zfs-lwb* and *zfs-pmem*, we configure the corresponding ZIL kind. For *zfs-sync_disabled*, we configure the ZIL-LWB ZIL kind but set the `sync=disabled` property. All variants use a `recordsize=4k` and `compression=off` on the dataset since most of our benchmarks are synchronous workloads with small write sizes. For the *zil-pmem* variant, we configure PRB with `ncommitters=3`.

{xfs,ext4}__on__\$BDEV__dax_\$DAX Linux 5.9's *xfs* or *ext4* filesystems deployed on the *block device stack* `$BDEV`. The `$DAX` value (True or False) indicates whether the `dax` mount option was set. We always perform benchmark runs without `dax`, and perform an additional run if `$BDEV` is a DAX-capable device.

make sure
DAX is ex-
plained

The *xfs* and *ext4* stacks are parametrized by the *block device stack* `$BDEV` which we define as a hardware and software configuration that provides a (potentially virtual) block device. We define the following block device stacks:

devpmem The raw `fsdax` namespace block device in `devfs`, i.e., `/dev/pmem0`.

dm-writocache The *dm-writocache* Linux Device Mapper target synchronously persists writes to PMEM and performs asynchronous writes-back to an origin block device in the background. We configure *dm-writocache* with `/dev/pmem0` as a cache device and a *dm-stripe* of the 30 NVMe partitions as the origin data store. To induce some level write-back, we configure *dm-writocache*'s *low watermark* to 0% and *high watermark* to 1%. (See Section 2.1.2 for details.)

zvol-lwb-rs_{0,1} A ZVOL exposed by a zpool in the same config as *zfs-lwb*. We set `volblocksize=4k` on the ZVOL dataset. The *rs* variable controls the value of the `zvol_request_sync` tunable which, if enabled, processes block I/O requests (`struct bio`) synchronously instead of submitting them to a thread pool. (Refer to Section 6.3.5 for more details on ZVOLs.)

zvol-sync_disabled Like *zvol-lwb*, but with `sync=disabled`.

zvol-pmem-rs_{0,1}-byp_{0,1}-nc_3 Like *zvol-lwb*, but configured with the ZIL-PMEM ZIL kind. PRB's `ncommitters` tunable is always set to 3. The *byp* variable controls whether the ITXG bypass is enabled (0 disabled, 1 enabled).

The following items are examples for storage stacks used in the evaluation:

ext4__on__zvol-pmem-rs_1-byp_0-nc_3__dax_False The Ext4 filesystem deployed on a ZVOL on a zpool with the ZIL-PMEM ZIL kind. The zpool is configured with `zvol_request_sync=1`, disabled ITXG bypass, and `ncommitters=3`. The ext4 file system is mounted without the `dax` mount option.

xfs__on__devpmem__dax_True The XFS filesystem, deployed directly on the PMEM block device (`/dev/pmem0`). It is mounted with the `dax` option.

zfs-pmem ZFS with the ZIL-PMEM ZIL kind and `ncommitters=3`. Note that `zvol_request_sync` and the ITXG bypass are only relevant for ZVOLs and therefore are not part of the configuration name.

We compare the performance of the storage stacks with the set of benchmarks described below. To facilitate the comparison, all benchmarks have the following properties:

Scaling factor The scaling factor increases the degree of concurrent synchronous I/O operations issued by the benchmark. We run each benchmark with three scaling factor values: 1, 4, and 8.

Result Metric Each benchmark reports a single result metric per run. This allows us to compare performance across different storage stacks and scaling factor values. For all benchmarks in this evaluation, a numerically greater result metric is better.

What follows is a description of the benchmarks:

fio-growing The 4k synchronous random write workload that motivated our work, but with all worker thread files within the same filesystem. The scaling factor maps to the `numjobs` parameter which controls the number of worker threads. The working set (amount of modified data) grows with the scaling factor because each additional thread adds a private file with a constant size of 100 MiB. We report mean IOPS as the result metric.

fio-fixed The 4k synchronous random write workload, with all files on a single filesystem, but a fixed working set size, independent of `numjobs`. We achieve the fixed working set size by setting the size parameter to $\frac{1 \text{ GiB}}{\text{numjobs}}$. We report mean IOPS as the result metric.

filebench varmail Filebench is a benchmarking engine that executes arbitrary workloads described in a domain-specific language. The predefined *varmail* workload “emulates I/O activity of a simple mail server that stores each e-mail in a separate file (`/var/mail/ server`). The workload consists of a multi-threaded set of create-append-sync, read-append-sync, read and delete operations in a single directory. [...] The workload generated is somewhat similar to Postmark but multi-threaded” []. We use filebench in version 1.5-alpha1-33-g22620e6 with the upstream workload definitions. The scaling factor maps to the workload’s `$nthreads` variable which determines the number threads that execute the opera-

tions described above. We use a runtime of 20 seconds per configuration and report filebench's *ops per second* value as the result metric.

filebench oltp The filebench *oltp* workload emulates database workloads. It "performs file system operations using Oracle 9i I/O model. It tests the performance of small random reads and writes, and is sensitive to the latency of moderate size (128k+) synchronous writes to the log file. By default launches 200 reader processes, 10 processes for asynchronous writing, and a log writer." []. We leave all parameters at the default setting except for the `$ndbwriters` parameter for which we use the scaling factor's value instead of its default value of ten. It is our understanding that scaling the number of log writer threads to a number greater than one would not be feasible with the real DBMS and thus be unrealistic. We use a runtime of 20 seconds and report filebench's *ops per second* value as the result metric.

MariaDB/sysbench We deploy the MariaDB 10.5.9 Docker image in its default configuration (InnoDB storage engine) with the `/var/lib/mysql` directory bind-mounted to a sub-directory within the storage stack's mountpoint. We apply the *sysbench* benchmark's *oltp_insert* workload with its default parameters for ten seconds. The workload spawns a number of threads that inserts rows with random data into one or more tables. We use the default setting (a single table) and map the scaling factor to the number of threads that perform the insert queries. Each thread uses a private, long-lived connection to the MariaDB server. We use Docker's `--net=host` parameter when starting the MariaDB container to allow *sysbench* to connect via loopback TCP, avoiding the overhead of Docker's user-space proxy. The result metric is the number of transactions per second (tps) reported by *sysbench*.

Redis-SET Redis is a popular in-memory key value store. For persistence, it provides two mechanisms that are recommended to be used in combination. First, the system periodically persists a snapshot of the Redis database (RDB). This process happens in the background in a forked child process. Second, Redis features a logical write-ahead log (*append-only file*, AOF) that is extended for every mutating operation and replayed after a crash. Redis supports three different behaviors for ensuring durability of the AOF, configurable through the `appendfsync` configuration variable. A value of `no` does not perform any fsyncs, `everysec` performs an fsync every second (default), and `always` performs an fsync operation for operation logged to AOF. []

We deploy Redis 6.2 (built from source). We configure RDB writeback to happen every second, regardless of the number of changes, We enable

AOF and configure `appendfsync=always` which effectively turns our Redis deployment into a durable key-value store.

For benchmarking, we use Redis's own *redis-benchmark* tool []. The scaling factor maps to the `--threads` and `-c` parameters which control the number of parallel clients. We use the benchmark's *SET* workload to perform 10^6 *SET* operations with random keys and the default value size (3 bytes). We configure a key space size of 10^6 to reduce contention. We leave the values for pipelining and connection keepalive at their defaults (no pipelining, keepalive enabled). The result metric is the requests-per-second (rps) value reported by *redis-benchmark*.

We choose the number of 10^6 *SET* operations because it results in runtimes of ten seconds or more for most combinations of scaling factor and storage stacks. The only exception is *zfs-sync-disabled* for which the shortest runtime is 8 s.

RocksDB-fillsync RocksDB is a popular key-value store developed by Facebook that is optimized for fast storage devices. It can be used directly by applications as an embedded database but is also the basis for the *MyRocks* storage engine for MySQL []. RocksDB uses log-structured merge trees for long-term storage which are written and rewritten in large units that are prepared in DRAM (memtables). Operations that request synchronous semantics using the `WriteOptions.sync` flag are logged to a write-ahead log file. [;]

We measure the impact of ZIL-PMEM on RocksDB WAL performance with the *fillsync* benchmark that is part of RocksDB's *db_bench* tool. *Fillsync* performs a fixed number of synchronous *Put* operations with random keys. We set the number of operations to 400k and map the scaling factor to the number of concurrently *Putting* threads. The result metric is the *ops-per-sec* (operations per second) value reported by *db_bench*.

We choose the number of 400k operations because it results in at least ten seconds of runtime for 90% of configurations. The configurations that achieve less than ten seconds of runtime are:

- All *zfs-pmem* configurations at scaling factor 1 (~4 s runtime)
- *zfs-sync-disabled* for scaling factors 1, 4, 8 (~2, 6.8 and 8 s rt.)
- *xfs-on-devpmem* with `dax` mount option at scaling factor 1 (9.1 s rt.)

7.3.4 ZIL-PMEM vs. ZIL-LWB

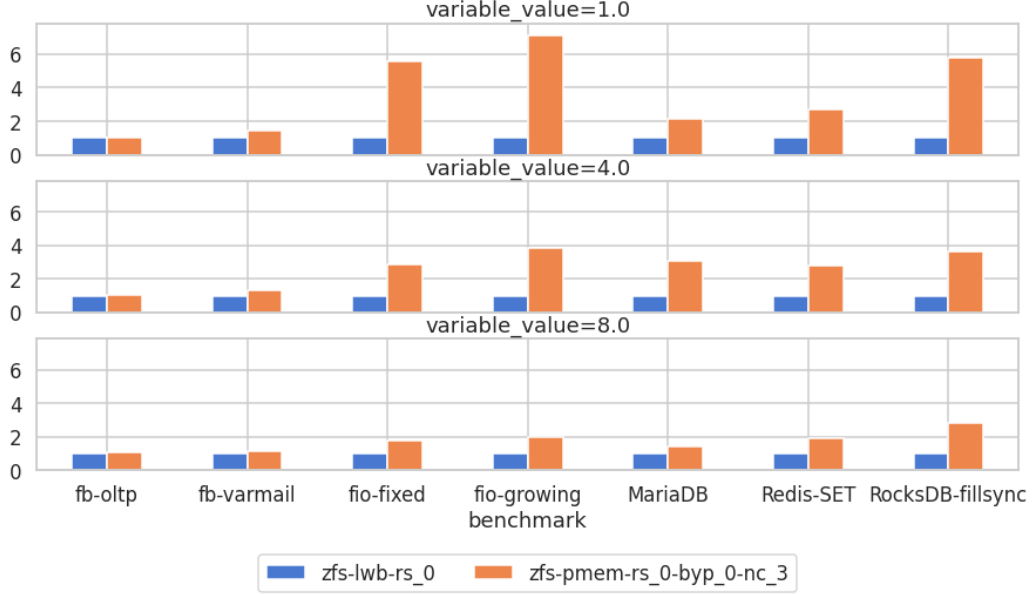


Figure 7.6: App Benchmarks: ZIL-PMEM vs. ZIL-LWB with ZIL-LWB as baseline.

We first compare the speedup of ZIL-PMEM over ZIL-LWB. ZIL-PMEM outperforms ZIL-PMEM in all benchmarks, but the speedup varies significantly for the different workloads and scaling factors. The highest speedups are achieved at scaling factor 1 by the fio workloads (7.1x, 5.3x) and *RocksDB-fillsync* (100k ops-per-sec, a speedup of 5.8). *Redis-SET* achieves shows a speedup of 2.68 (30k rps) and MariaDB achieves 10k tps which is a speedup of 2.14. For scaling factor 4, the fio and RocksDB's speedup declines whereas Redis's grows to 2.77 and MariaDB's grows to 3.08. For scaling factor 8, all workloads show a reduction in speedup (RocksDB 2.8x, Redis 1.9x, MariaDB 1.4x).

The *filebench-oltp* workload shows no relevant speedup or slowdown for all scaling factor values. We observed during benchmark execution that the 200 reader processes always cause 100% CPU utilization.

The *filebench-varmail* workload only shows a speedup of 1.4 at scaling factor 1 which shrinks to 1.1 at scaling factor 8. A possible explanation for these meager results is the larger data volume (16k appends) compared to the fio workloads (4k writes). In combination with the large amount of metadata (file creation and deletion), ZIL-LWB might also be benefitting from the *commit ITXs*' pipelining. In any way, *filebench-varmail* shows that the large speedups for small writes

cannot be generalized to all filesystem workloads. Future work should investigate ZIL-PMEM's behavior during this benchmark in detail and determine how beneficial pipelining akin to *commit ITXs* could be.

7.3.5 ZIL-PMEM vs. XFS and Ext4 on PMEM

Our next set of results is the comparison between ZFS and the Linux filesystems XFS and Ext4. We include both ZIL kinds (*lwb*, *pmem*) as well as `sync=disabled` to demonstrate the upper bound for any ZIL kind. Both XFS and Ext4 are DAX-aware and thus we include configurations with and without the *dax* mount option. Note that that the ZFS configurations only use PMEM for the ZIL but NVMe devices for permanent storage whereas the Linux filesystem configurations are PMEM-only.

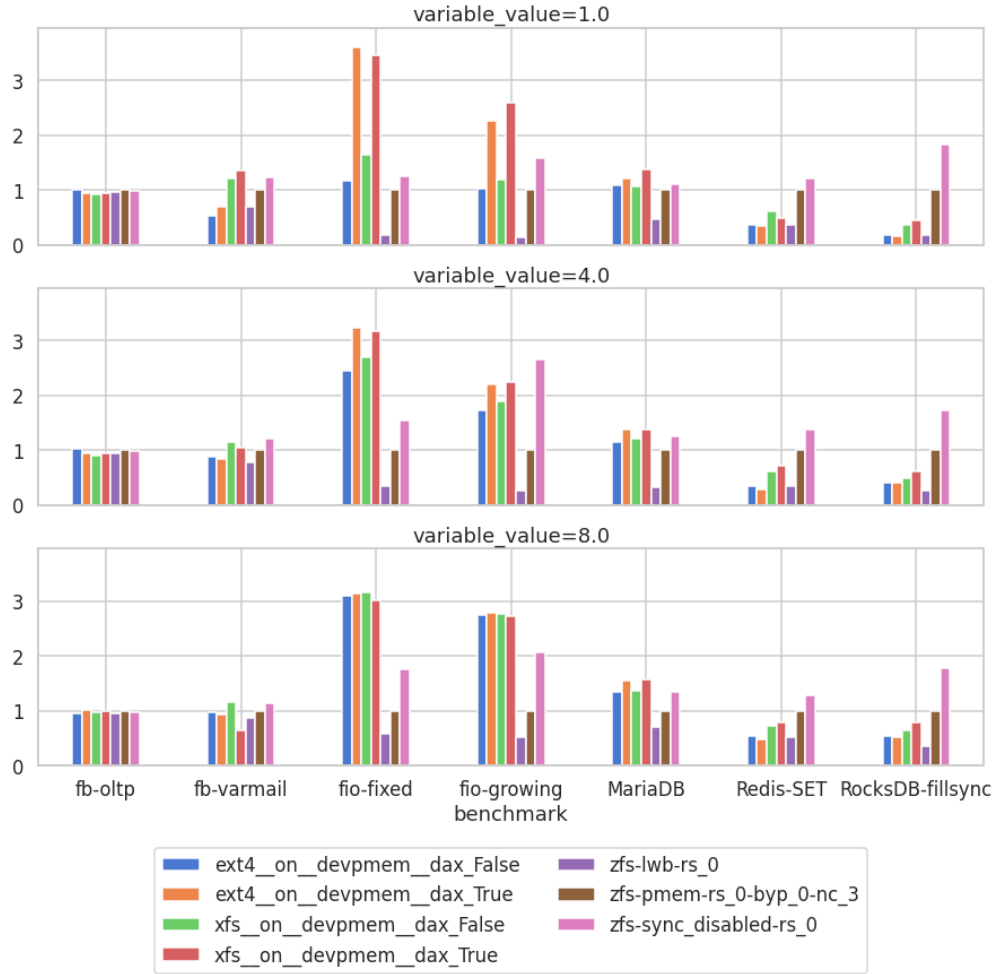


Figure 7.7: Speedups of ZFS with ZIL-LWB and sync=disabled as well as ext4 and xfs with and without the *dax* mount option over ZIL-PMEM as baseline.

Our primary takeaway from this comparison is that, in contrast to ZIL-LWB, ZIL-PMEM is competitive with both ext4 and xfs at scaling factor 1 in most benchmarks. The fio workloads are the exception, where the Linux filesystems achieve a speedup of greater than 3 if the *dax* mount option is enabled.

ZIL-PMEM performs significantly better than all other configurations (excluding sync=disabled) for *Redis-SET* and *RocksDB-fillsync* at all scaling factors. The write-ahead log files written by these workloads are append-only and contain the encoded representation of logical changes. The mutating requests issued by our

benchmarks have small payloads and result in small `write+fsync` system calls.¹ ZIL-PMEM handles this workload more efficiently than the Linux filesystems because the ZIL's data journaling causes less write amplification than the Linux filesystems whose latency is likely dominated by page-granular write-back on every `fsync`.

In the *MariaDB* benchmark, the DAX-aware Linux filesystem configurations achieve up to 1.6x more transactions per second than ZIL-PMEM. The `sync=disabled` workload is on-par with the non-DAX-aware configurations, at speedups over ZIL-PMEM of ca. 1.1, 1.2 and 1.3 at scaling factors 1, 4, and 8. Manual inspection of the CPU utilization shows a significant amount of context switching but no clearly identifiable CPU bottleneck. The database server process performs 1024 `write + fdatasync` system calls to the InnoDB redo log file from a single thread. However, since `sync=disabled` does not perform better than any of the systems that actually write to PMEM, we can also rule out a PMEM bandwidth as a possible bottleneck.

Regarding the *filebench-oltp* benchmark, no configuration performs exceptionally better than any other, including `sync=disabled`. We observe high CPU utilization in all benchmark configurations, caused by the *hog* stage of the 200 processes that simulate reader threads. Thus, CPU is most likely the bottleneck in all configurations.

need to state that we do not have a final conclusion here?

For *filebench-varmail*, the highest speedup is 1.35 at scaling factor 1, achieved by XFS in DAX-aware mode, followed by ZFS with `sync=disabled` with 1.22x. For higher scaling factors, the speedups are slightly lower. The fact that `sync=disabled` achieves a lower speedup than XFS can be taken as an indicator that *filebench-varmail* is not bound by PMEM performance, contradicting our speculation in the previous section.

7.3.6 ZIL-PMEM vs. XFS and Ext4 on Dm-writecache

A fairer comparison is thus the deployment of the Linux filesystems on top of *dm-writecache* with a *dm-stripe* of the NVMe partitions as the origin block device. The comparison between ZIL-PMEM with NVMe storage and Linux filesystems on pure PMEM in the previous section is insightful but not a level comparison. The ZFS configuration only uses PMEM as short-term storage to bridge the gap

¹The RocksDB documentation and comments in the source code suggest that the WAL files are written in 32k blocks that are zero-padded if necessary [`db/log_writer.h`]. We observed the contrary behavior with `strace`: each *Put* operation with `WriteOptions.sync` results in a `write` and `fdatasync` system call. The `write` system call's data size is proportional to the sum of the *Put* operation's key and value sizes.

between the background txg syncs to the striped pool of NVMe partitions. In contrast, the Linux filesystem configurations store all their data on PMEM. Essentially, this

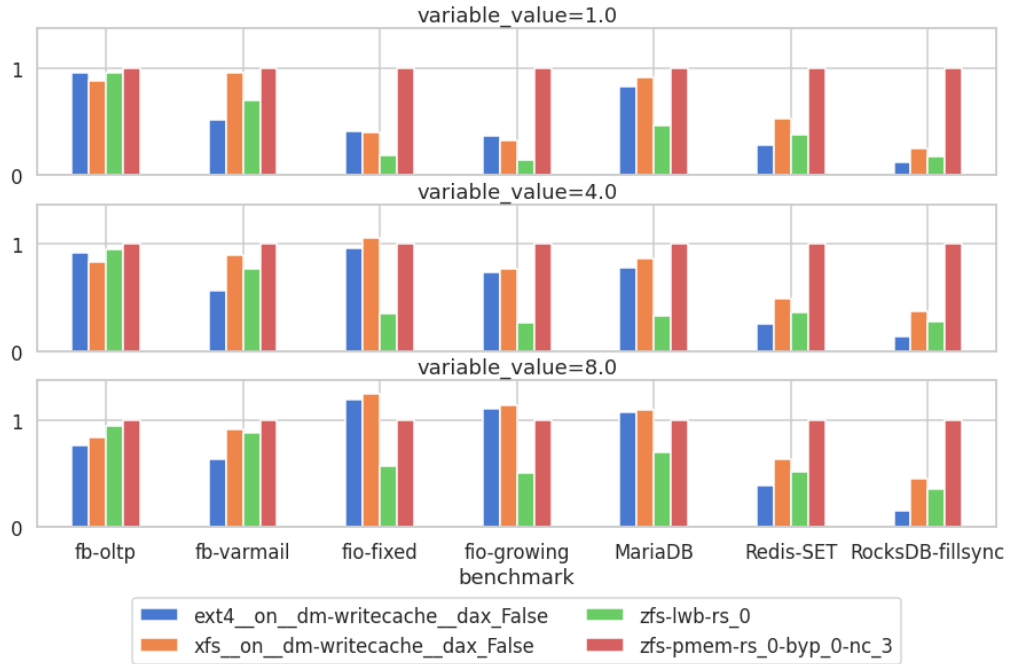


Figure 7.8: Comparison of Linux filesystems on dm-writecache with ZIL-PMEM as a baseline.

ZIL-PMEM's performance is within $\pm 30\%$ of XFS on *dm-writecache* at all scaling factors in all benchmarks except *Redis-SET* and *rocksdb-fillsync* where ZIL-PMEM performs significantly better. The 30% spread decreases slightly with higher scaling factor but does not change the overall impression. Our explanation for ZIL-PMEM's significantly better performance in the *Redis-SET* and *rocksdb-fillsync* workloads is (again) write amplification: due to ZFS's integrated design (filesystem and volume manager), the ZIL is able to write logical changes directly to PMEM with little write amplification. In contrast, dm-writecache hides PMEM behind the block device interface, prohibiting direct access to PMEM (i.e., the *dax* mount option is not available).

Ext4 performs significantly worse than ZIL-PMEM (and XFS!) in some benchmarks, with only 45%, 28%, and 14% of ZIL-PMEM's performance in *filebench-varmail*, *Redis-SET*, and *rocksdb-fillsync* at scaling factor 1. Write amplification might be a possible explanation as well.

better word?

We have manually observed the system while benchmarking the *dm-writecache* stack to ensure that the comparison with ZFS is actually fair. We found that, despite the very aggressive write-back configuration (high watermark = 1, low watermark = 0), the system performed very little write-back to the NVMe drives. The two *fio* workloads were the exception but the write load towards the NVMe devices was well below their maximum capacity. Thus, NVMe was not a bottleneck in the benchmark. However, we found that *dm-writecache* has severe multi-core scalability problems. For example, XFS on raw PMEM without the *dax* mount option achieves 150k IOPS at scaling factor 1 and 409k IOPS at scaling factor 4. In contrast, on *dm-writecache*, it is 112k IOPS at scaling factor 1 but only 165k IOPS at scaling factor 4. Using the *perf* tool, we observed severe contention at a lock that serializes access to the entire *dm-writecache* instance's state. In private communication, *dm-writecache*'s maintainer Mikulas Patocka stated that "the purpose of *dm-writecache* is to decrease commit latency, not to increase throughput. There is not much that can be done with the lock contention". ZIL-PMEM has similar problems (125k at factor 1, 217k at factor 4) albeit due to the sequential structure of ZIL chain. In the next section, we show how a more relaxed log structure such as the ITXG bypass can improve the situation.

waiting for permission to quote him publicly

7.3.7 Impact on ZVOL Performance & ITXG Bypass

In this section, we examine the impact of ZIL-PMEM on ZVOL performance. We compare the performance of XFS deployed on a ZVOL in a zpool with varying ZIL kinds, *zvol_request_sync* (*rs_{0,1}*), and ITXG bypass setting (*byp_{0,1}*). The results are relative to *zvol-pmem-rs_0-byp_0-nc_3* as the baseline which is the standard configuration for ZIL-PMEM pools. Note that deploying a filesystem on a ZVOL generally makes little sense in practice because ZFS filesystems provide the same basic service with more additional features at less overhead. However, ZVOLs are a popular choice for storage virtualization where ZVOLs are used as virtual hard disks, either on the same host or via a SAN (e.g., iSCSI, Fibre Channel). To simplify our evaluation, we avoid the overhead of a hypervisor or SAN and instead instantiate XFS directly on top of the ZVOL, mount it, and execute the benchmarks.

is this a good transition?

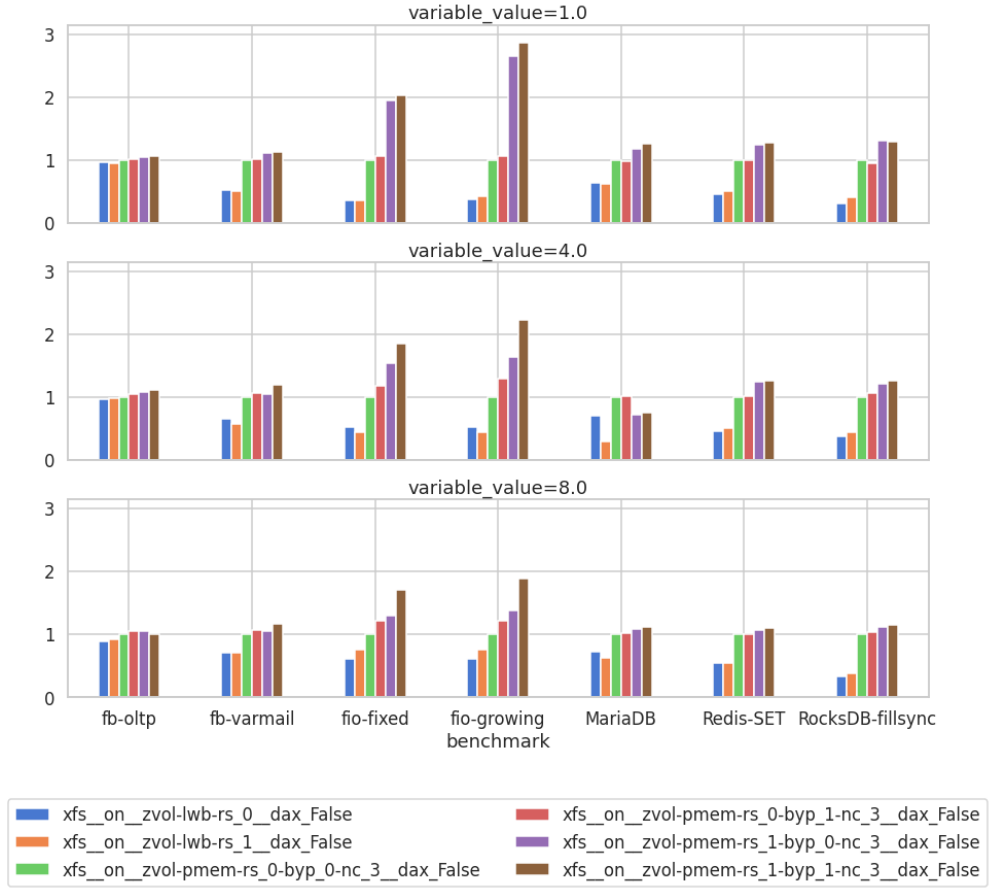


Figure 7.9: Comparison of XFS on ZVOLs on different zpool configurations.

ZIL-PMEM in its default configuration delivers a significant speedup over ZIL-LWB. At scaling factor 1, it is approximately 2x for all workloads except *MariaDB* ($\frac{1}{0.68} = 1.47$) and *filebench-oltp* ($\frac{1}{0.96} = 1.04$). At higher scaling factors, the speedup declines for all benchmarks albeit less so with *Redis-SET* and *RocksDB-fillsync* than with the other workloads.

The effect of `zvol_request_sync=1` is ambiguous and not very significant for ZIL-LWB but very beneficial for ZIL-PMEM in the *fio* workloads at scaling factor 1 with a 2x and 2.6x speedups over the standard configuration, respectively. *MariaDB*, *Redis-Set* and *RocksDB-fillsync* also benefit with speedups of 1.18x, 1.23x, and 1.30x. However, for scaling factor 4, *MariaDB* performs substantially worse if `zvol_request_sync` is set (30k tps vs. 20k tps). We also conducted the benchmarks with Ext4 instead of XFS (not shown in the plot). Ext4 exhibited worse

performance in almost all configurations if `zvol_request_sync` was set, with the exception of the *fio* workloads.

The ITXG bypass only has marginal effects with and without `zvol_request_sync`, except for the *fio* workloads at scaling factors 1 and 4. For example, *fio-growing* at scaling factor 8 shows an increase in IOPS from 84k to 102k (21%) for `zvol_request_sync=0` and from 117k to 159k (35%) with `zvol_request_sync=1`. For Ext4, *filebench-varmail* also shows a 20% improvement with enabled ITXG bypass and `zvol_request_sync=0`. Activation of the ITXG bypass does not appear to have a negative impact for XFS in the remaining workloads. For Ext4, we observed some performance degradations in setups whether both ITXG bypass and `zvol_request_sync` were enabled.

In summary, ZIL-PMEM provides a significant performance advantage for ZVOLs but is less effective than with ZFS filesystems (see Sections 7.3.2 and 7.3.4). This is particularly noticable for the *fio* workloads and *RocksDB-fillsync*. Their speedup with ZIL-PMEM on XFS+ZVOL is only half the speedup of what we observed for ZFS at scaling factor 1. Write amplification cannot be responsible for this behavior because only *RocksDB-fillsync* causes write amplification. The most likely cause is latency overhead added by the filesystem that affects all workloads equally.

7.3.8 CPU-Efficient Handling Of PMEM Bandwidth Limits

PRB's *commit slot* mechanism limits the amount of parallel writers to PMEM through the `ncommitters` variable. The goal is to avoid wasteful on-CPU stall cycles which inevitably occur if the aggregate write bandwidth to PMEM exceeds the hardware capacity. (See Sections 4.1.1 and 5.13 for a more detailed explanation.)

To determine the mechanism's effectiveness, we use our 4k synchronous random write workload with separate ZFS filesystems per *fio* thread. We add low-overhead per-CPU counters to ZIL-PMEM to sum up the total amount of time that is spent writing PMEM (T_{pmem}), as well as the number of write operations performed during the benchmark (N_{ops}). We then compute the *average PMEM write time per IOP* $T_{pmem_{iop}} = \frac{T_{pmem}}{N_{ops}} [\frac{s}{op}]$. The most efficient value for `ncommitters` for a given system then depends on the workload and efficiency goals. In general, `ncommitters` should be the minimal value where $T_{pmem_{iop}}$ is tolerable but the system's primary performance metric is not impacted.

We compute $T_{pmem_{iop}}$ for 1 to 18 numjobs, different values of `ncommitters` and two different PMEM configurations. The first PMEM configuration is a single, non-interleaved Optane DIMM. The second configuration are four interleaved Optane

DIMMs. Regardless of how interleaving is configured, we create a 40 GiB-sized *fsdax* namespace on the resulting region and use it as the PMEM SLOG for ZIL-PMEM. We use a machine with higher core count for the benchmarks. The system configuration is as follows:

System Supermicro SYS-1029U-TRT
Mainboard Supermicro X11DPU, Version 1.10
CPU 2 x Intel(R) Xeon(R) Gold 5220 CPU 2.20GHz
 18 cores per socket, 2-way SMT per core
DRAM 12 x 32GiB SK HYNIX DDR4-2666, HMA84GR7CJR4N-VK
PMEM 8 x Intel Optane DC Persistent Memory, 128 GB, (NMA1XXD128GPS),
 4 per socket.
NVMe System rootfs only, no dedicated NVMe hardware.
Kernel Linux Kernel, Fedora 5.11.15-100.fc32.x86_64
Userland Fedora 32
fio fio-3.21

As with our main evaluation system (described in Section 3.1), we leave SMT enabled. We disable the second socket in software using the `isolcpus=18-35,54-71` kernel command line parameter. Since the system does not have NVMe devices available, we provision the disabled socket’s Optane DIMMs in `AppDirectNotInterleaved` mode with one PMEM region and a single *fsdax* namespace per DIMM. We use these namespaces as the *zpool*’s main (non-SLOG) *vdevs*.

Figure 7.10 visualizes the results of this experiment. For the non-interleaved configuration, we observe that the system’s peak IOPS is 400k, first achieved with `ncommitters=3`. For `ncommitters=3`, IOPS quickly decline to 263k IOPS for higher *numjobs*. With `ncommitters=8`, the system is able to sustain the 400k IOPS for the for *numjobs* $\in 7 \dots 13$. The price is significantly more on-CPU PMEM time per IOP: looking at *numjobs* = 8, the $T_{pmem_{iops}}$ is 3300, 4300, 6500 ns for `ncommitters=3`, 4 and 8. And at *numjobs* = 12, $T_{pmem_{iops}}$ climbs to 11600, 20600 ns for `ncommitters=8` and 12 whereas `ncommitters=3` is successfully limited to 2700 ns of PMEM write time per IOP. The `ncommitters=12` and 24 configurations do not achieve higher peak IOPS than `ncommitters=8` but decline to lower values for high *numjobs*, e.g., only 237k IOPS for `ncommitters=24`, *numjobs* = 18. This performance decline at high degrees of concurrency is a property of the Intel Optane PMEM hardware [Yan+20]. The limitation to `ncommitters=8` mitigates this effect (331k vs. 237k IOPS at *numjobs* = 18).

The four-way interleaved configuration exhibits a very different behavior. We achieve the highest IOPS for the configuration where the *commit slot* mechanism has no effect, i.e., `ncommitters=24`, *numjobs* = 18 with 900k IOPS. Given that the IOPS curve has not reached a plateau at this point, a higher value might be

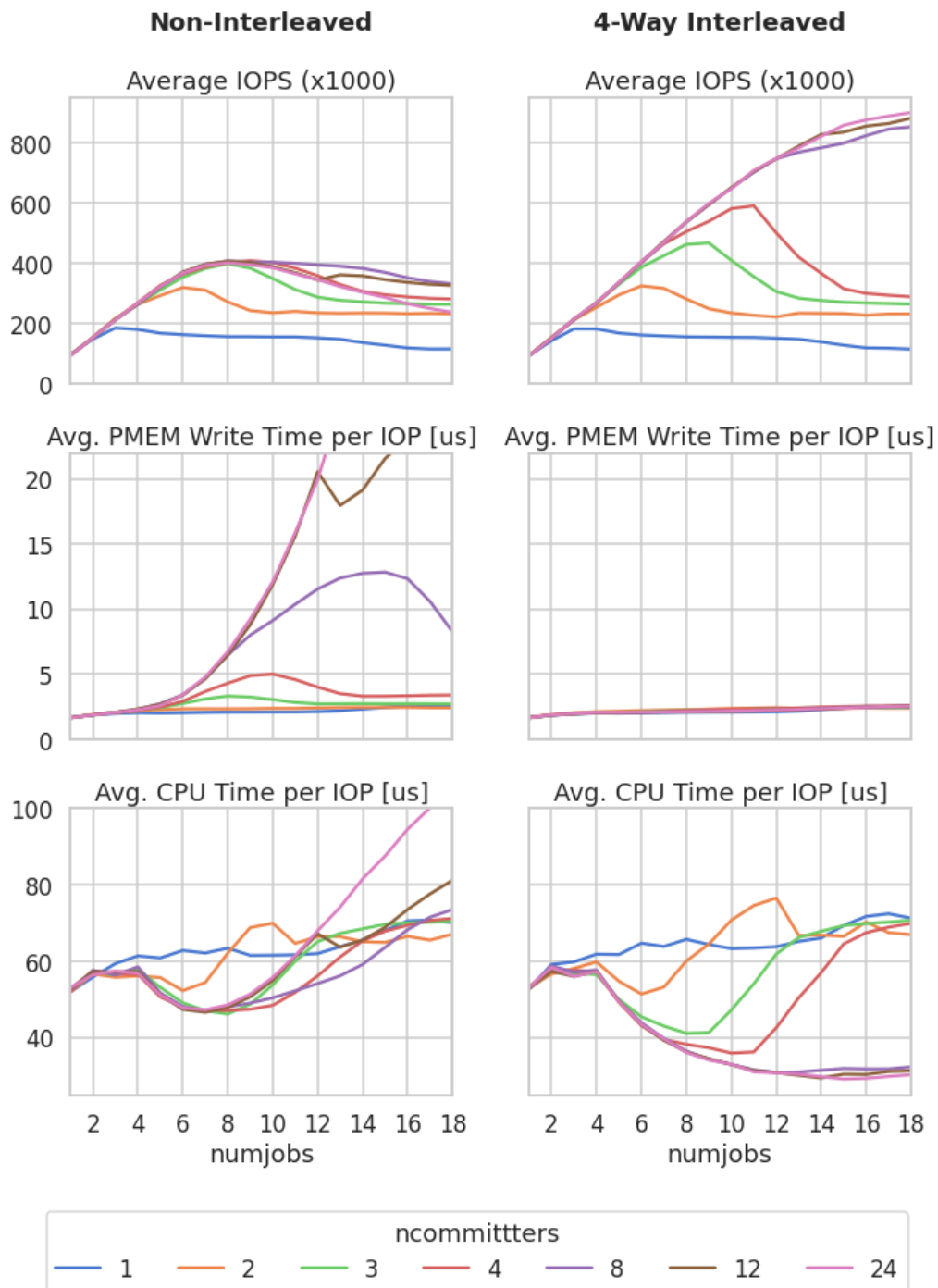


Figure 7.10: Comparison of performance and PMEM write time per IOPS for different values of `ncommitters` in two PMEM configurations (left non-interleaved, right four-way interleaved).

possible with higher $numjobs$. $T_{pmem_{iops}}$ is the same for all $ncommitters$ values at a given $numjobs$ value (± 100 ns). It is 1660 ns for $numjobs = 1$ and grows slightly unevenly towards 2560 ns at $numjobs = 24$. This is an increase of merely 54% which is likely to be acceptable in any setup, given the 9.6x increase in IOPS (93k to 900k) that is possible with $ncommitters=24$.

To determine the *commit slot* abstraction's impacts performance in $ncommitters > numjobs$ configurations, we add additional instrumentation that measures the average latency of commit slot acquisition and release. Figure 7.11 visualizes the results for $ncommitters=24$. The absolute overhead is approximately the same for both PMEM configurations. It starts at 100 ns for $numjobs = 1$, climbs to 240 ns at $numjobs = 4$, and then scales linearly to 385 ns at $numjobs = 18$. In the interleaved configuration, for any $numjobs > 4$, this corresponds to approximately 11–13% of the overall write latency per log entry. Without *commit slots*, ztl_commit could thus be $\frac{1}{1-0.13} = 14.9\%$ faster. However, due to the overhead of the other ZFS components, the contribution to overall IOP latency is only 2%, which would constitute a rather meager improvement.

Think this paragraph through again, I think it's correct, but should be double-checked.

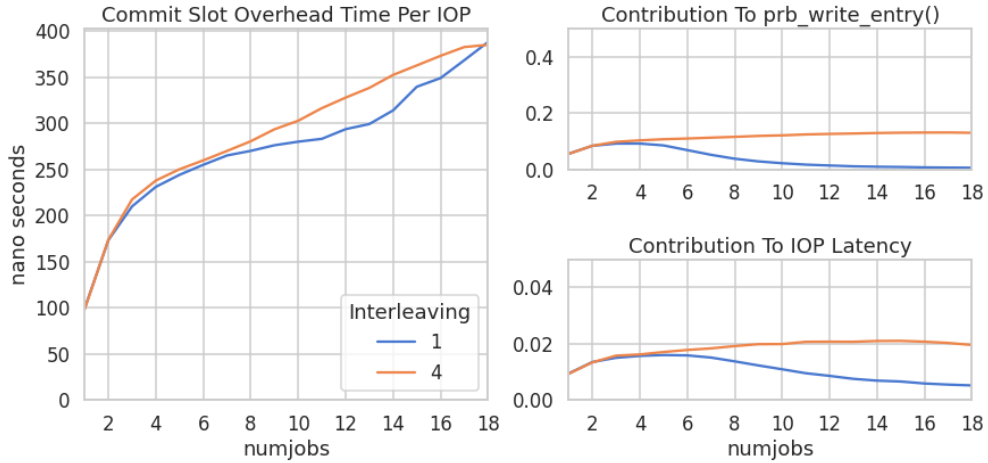


Figure 7.11: Overhead of the *commit slot* for $24 = ncommitters > numjobs = 18$. Note the different y-axis scales. Also note that the relative latency contribution for the non-interleaved configuration only shrinks because the PMEM write time (not shown) grows significantly.

We draw the following conclusion from the observations made above:

- The *commit slot* mechanism successfully limits PMEM write time.

- The *commit slots* mechanism's multi-core scalability is acceptable and contributes little relative overhead to each IOP compared to other ZFS components.
- Regarding the non-interleaved configuration (single Optane DIMM):
 - Our 4k synchronous write workload is able to achieve peak performance with 400k IOPS as early as `ncommitters=3`.
 - Setting `ncommitters=8` extends the peak to a plateau with regards to *numjobs*, but consumes $\frac{6500 \text{ ns}}{3300 \text{ ns}} = 1.96$ times more (on-CPU!) PMEM write time per IOP.
 - Note that we use `ncommitters=3` as the default configuration for ZIL-PMEM in all of the previous benchmarks. Our rationale is that an in-kernel filesystem should not waste CPU time on PMEM I/O because it is a system-wide OS service whose overall priority in the system is not known ahead of time. However, if the system's primary role is that of an NFS/SMB share, or that of a SCSI target (ZVOLS), it may be appropriate to allow more on-CPU waiting.
- Regarding the four-way interleaved configuration (4 Optane DIMMs):
 - The *commit slot* mechanism is not useful for the observed range of *numjobs* values.
 - The reason is that ZIL-PMEM cannot saturate the Optane DIMMs' write bandwidth due to per-IOP overhead added by other components of ZFS.
 - However, with `ncommitters > numjobs`, the overhead per IOP is negligible.

7.4 Discussion

Storage stacks that we investigated but showed subpar perf:

- xfs with separate log
- ext4 with separate log

Future work:

- proof our write amplification claims (measure number of bytes of PMEM written as well as pmem write time, for all configurations. Needs instrumentation of pmem block device emulation.)
- ext4/xfs data journaling mode (ext4 has it, xfs not sure)
- xfs pmem-aware log (Christoph Hellwig has patches...)

- Scalability comparison of ZIL-LWB and ZIL-PMEM within the same dataset. Expectation: ZIL-PMEM's mutex eventually becomes the bottleneck, but likely still better than ZIL-LWB due to massive latency headroom.

Chapter 8

Summary

8.1 Future Work

Evaluation Of Replay Performance

Evaluation of Chunk Size

Appendix

Bibliography

- [] *[Ndctl PATCH v2 3/6] Ndctl: Add an Inject-Error Command - Linux-Nvdim - Ml01.01.Org*. URL: <https://lists.syncevolution.org/hyperkitty/list/linux-nvdim@lists.01.org/message/UAEXUT4RBYFUFTXIFPMH6OZM6FWYD7A4/> (visited on 04/13/2021).
- [] *Fast Commits for Ext4 [LWN.Net]*. URL: <https://lwn.net/SubscriberLink/842385/ea43ae3921000c72/> (visited on 01/16/2021).
- [] *Filebench GitHub Wiki / Predefined Personalities*. GitHub. URL: <https://github.com/filebench/filebench/wiki/Predefined-personalities> (visited on 03/30/2021).
- [] *FlushWAL; Less Fwrite, Faster Writes*. RocksDB. URL: <http://rocksdb.org/blog/2017/08/25/flushwal.html> (visited on 10/03/2020).
- [] *How Fast Is Redis? – Redis*. URL: <https://redis.io/topics/benchmarks> (visited on 05/15/2021).
- [] *Linux 5.0 Compat: SIMD Compatibility · Openzfs/Zfs@e5db313*. GitHub. URL: <https://github.com/openzfs/zfs/commit/e5db31349484e5e859c7a942eb15b98d68ce5b4d> (visited on 05/04/2021).
- [] *Metadata Allocation Classes by Don-Brady · Pull Request #5182 · Openzfs/Zfs*. GitHub. URL: <https://github.com/openzfs/zfs/pull/5182> (visited on 04/13/2021).
- [] *MyRocks | A RocksDB Storage Engine with MySQL*. MyRocks. URL: <http://myrocks.io/> (visited on 05/16/2021).
- [] *Ndctl-Inject-Error Man Page - Ndctl - General Commands*. URL: <https://www.mankier.com/1/ndctl-inject-error> (visited on 04/13/2021).
- [] *Redis Persistence – Redis*. URL: <https://redis.io/topics/persistence> (visited on 05/15/2021).
- [] *RocksDB GitHub Wiki: WAL Performance*. GitHub. URL: <https://github.com/facebook/rocksdb/wiki/WAL-Performance> (visited on 05/15/2021).
- [] *Valgrind*. URL: <https://www.valgrind.org/> (visited on 04/13/2021).

- [] *Writecache Target — The Linux Kernel Documentation*. URL: <https://www.kernel.org/doc/html/latest/admin-guide/device-mapper/writecache.html> (visited on 04/13/2021).
- [] *Zinject OpenZFS Man Page*. GitHub. URL: <https://github.com/openzfs/zfs/blob/65c7cc49bfcf49d38fc84552a17d7e8a3268e58e/man/man8/zinject.8> (visited on 04/13/2021).
- [Bon+03] Jeff Bonwick et al. “The Zettabyte File System.” In: *Proc. of the 2nd Usenix Conference on File and Storage Technologies*. Vol. 215. 2003.
- [Bor+16] James Bornholt et al. “Specifying and Checking File System Crash-Consistency Models.” In: *ACM SIGPLAN Notices* 51.4 (Mar. 25, 2016), pp. 83–98. ISSN: 0362-1340. DOI: 10.1145/2954679.2872406. URL: <https://doi.org/10.1145/2954679.2872406> (visited on 02/03/2021).
- [Con+09] Jeremy Condit et al. “Better I/O through Byte-Addressable, Persistent Memory.” In: *Proceedings of the ACM SIGOPS 22nd Symposium on Operating Systems Principles*. SOSP ’09. New York, NY, USA: Association for Computing Machinery, Oct. 11, 2009, pp. 133–146. ISBN: 978-1-60558-752-3. DOI: 10.1145/1629575.1629589. URL: <https://doi.org/10.1145/1629575.1629589> (visited on 02/04/2021).
- [DeW+84] David J. DeWitt et al. “Implementation Techniques for Main Memory Database Systems.” In: *Proceedings of the 1984 ACM SIGMOD International Conference on Management of Data*. 1984, pp. 1–8.
- [Dul+14] Subramanya R. Dulloor et al. “System Software for Persistent Memory.” In: *Proceedings of the Ninth European Conference on Computer Systems*. EuroSys ’14. New York, NY, USA: Association for Computing Machinery, Apr. 14, 2014, pp. 1–15. ISBN: 978-1-4503-2704-6. DOI: 10.1145/2592798.2592814. URL: <https://doi.org/10.1145/2592798.2592814> (visited on 02/04/2021).
- [Fan+11] Ru Fang et al. “High Performance Database Logging Using Storage Class Memory.” In: *2011 IEEE 27th International Conference on Data Engineering*. 2011 IEEE 27th International Conference on Data Engineering. Apr. 2011, pp. 1221–1231. DOI: 10.1109/ICDE.2011.5767918.
- [HLM94] Dave Hitz, James Lau, and Michael A. Malcolm. “File System Design for an NFS File Server Appliance.” In: *USENIX Winter*. Vol. 94. 1994.
- [Joh+10] Ryan Johnson et al. “Aether: A Scalable Approach to Logging.” In: *Proceedings of the VLDB Endowment* 3.1-2 (2010), pp. 681–692.
- [Kad+19] Rohan Kadekodi et al. “SplitFS: Reducing Software Overhead in File Systems for Persistent Memory.” In: *Proceedings of the 27th ACM Symposium on Operating Systems Principles*. SOSP ’19. New York, NY, USA: Association for Computing Machinery, Oct. 27, 2019,

- pp. 494–508. ISBN: 978-1-4503-6873-5. DOI: 10.1145/3341301.3359631. URL: <https://doi.org/10.1145/3341301.3359631> (visited on 10/01/2020).
- [Kwo+17] Youngjin Kwon et al. “Strata: A Cross Media File System.” In: *Proceedings of the 26th Symposium on Operating Systems Principles*. SOSP ’17. New York, NY, USA: Association for Computing Machinery, Oct. 14, 2017, pp. 460–477. ISBN: 978-1-4503-5085-3. DOI: 10.1145/3132747.3132770. URL: <https://doi.org/10.1145/3132747.3132770> (visited on 10/01/2020).
- [Lan+14] Philip Lantz et al. “Yat: A Validation Framework for Persistent Memory Software.” In: *2014 USENIX Annual Technical Conference (USENIX ATC 14)*. 2014, pp. 433–438.
- [LBN13] Eunji Lee, Hyokyung Bahn, and Sam H. Noh. “Unioning of the Buffer Cache and Journaling Layers with Non-Volatile Memory.” In: *Presented as Part of the 11th USENIX Conference on File and Storage Technologies (FAST 13)*. 2013, pp. 73–80.
- [Liu+19] Sihang Liu et al. “Pmtest: A Fast and Flexible Testing Framework for Persistent Memory Programs.” In: *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 2019, pp. 411–425.
- [Ope20] OpenZFS, director. *ZIL Performance Improvements for Fast Media by Saji Nair*. Oct. 12, 2020. URL: <https://www.youtube.com/watch?v=TnXwrigwF7I> (visited on 04/13/2021).
- [PAA05] Vijayan Prabhakaran, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. “Model-Based Failure Analysis of Journaling File Systems.” In: *2005 International Conference on Dependable Systems and Networks (DSN’05)*. IEEE, 2005, pp. 802–811.
- [Pel+13] Steven Pelley et al. “Storage Management in the NVRAM Era.” In: *Proceedings of the VLDB Endowment 7.2* (2013), pp. 121–132.
- [Pil+14] Thanumalayan Sankaranarayana Pillai et al. “All File Systems Are Not Created Equal: On the Complexity of Crafting Crash-Consistent Applications.” In: *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*. 2014, pp. 433–448.
- [PS17] Daejun Park and Dongkun Shin. “iJournaling: Fine-Grained Journaling for Improving the Latency of Fsync System Call.” In: *2017 USENIX Annual Technical Conference (USENIX ATC 17)*. 2017, pp. 787–798.
- [Sca20] Steve Scargall. “Persistent Memory Architecture.” In: *Programming Persistent Memory: A Comprehensive Guide for Developers*. Berkeley, CA: Apress, 2020, pp. 11–30. ISBN: 978-1-4842-4932-1. DOI: 10.1007

