

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 4. Mai 2021

Abstract

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	21
2.3 OpenZFS Primer	21
2.3.1 The ZIL API	21
3 Why ZIL-LWB Is Slow On PMEM	25
4 Design & Implementation	27
4.1 Project Scope	27
4.1.1 Requirements	27
4.1.2 Out Of Scope For The Thesis	29
4.1.3 Limitations	29
4.2 Overview	30
4.3 ZIL kinds	32
4.3.1 On-Disk State	33
4.3.2 Runtime State	34
4.3.3 Changing ZIL Kinds	35
4.3.4 ZIL-LWB Suspend & Resume	38

4.3.5	ZIL Traversal & ZDB	38
4.3.6	ZIL-LWB-Specific Callbacks	39
4.3.7	Considered Alternatives	40
4.3.8	Summary	42
4.4	PMEM-aware SPA & VDEV layer	42
5	PRB/HDL	45
5.1	OpenZFS Background	45
5.2	Analysis	48
5.3	Approach	49
5.4	HDL: Log Structure	50
5.5	HDL Replay: Basic Approach	51
5.6	HDL Replay: Dependency Tracking	52
5.7	HDL: Model For Data Corruption	58
5.8	HDL: Replay Crash-Consistency	59
5.9	PRB: DRAM Data Structure	60
5.10	PRB: PMEM Data Structure	64
5.11	PRB: Chunk Traversal	65
5.12	PRB: Write Path	68
5.12.1	Commit Slots	69
5.12.2	HDL-Scoped Metadata	72
5.12.3	Crash-Consistent Insert	72
5.13	API Design	75
5.13.1	PRB Setup	75
5.13.2	Replay	77
5.13.3	Writing Entries	78
6	ZIL-PMEM	81
7	ITXG Bypass For ZVOL	83
8	Evaluation	85
8.1	Correctness	85
8.1.1	PRB	85
8.1.2	ZIL-PMEM	85
8.2	Performance	85
8.2.1	Write Performance	85
8.2.2	Replay Performance	85
9	Related Work	87
	Appendix	89

<i>CONTENTS</i>	3
Bibliography	91

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 OpenZFS 2.0, 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.

review numbers in this paragraph

Our investigation of the latency distribution for the OpenZFS PMEM SLOG configuration shows that the ZIL implementation contributes 48 us in the single-threaded case and up to 95 us at eight threads. More specifically, up to 50 % of total syscall latency is spent waiting for the *ZIO pipeline*, excluding the time spent in the PMEM block device driver.

run this again without ebpf

ZIO is ZFS's abstraction for performing IO to the pool's storage devices. The pipeline-oriented design has proven extremely flexible and extensible over ZFS's lifetime. Many of ZFS's distinguishing features are implemented in ZIO, including physical block allocation, transparent data and metadata checksumming, compression, deduplication and encryption as well as redundancy mechanisms such as raidz. CPU-intensive processing steps are parallelized to all of the system's CPUs using *taskqs*.

It is our impression that ZIO is biased towards throughput, not latency. For example, ZIO's main consumer is *txg sync* which issues the updates to the on-disk merkle tree as a large tree of dependent ZIOs. Other consumers such as periodic data integrity checking (*scrubbing*) are also throughput-oriented. In contrast, the ZIL is the only component that is latency-oriented. ZIO's latency overhead with fast block devices such as NVMe drives is a well known problem among ZFS developers.

Given

ref, gh issue, cite private conversation with brian?

- the abysmal out-of-the-box performance of the PMEM block device as a SLOG, and
- our findings on the ZIL's and ZIO pipeline's significant latency overhead, as well as the fact that
- grouping log records into log-write *blocks* is unnecessary with byte-addressable PMEM,

too harsh?

right prep?

we believe that the current ZIL design is fundamentally unfit to reap the performance available with PMEM. We present **ZIL-PMEM**, a new implementation of the ZIL that bypasses the ZIO pipeline completely and persists log entries directly to 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 8 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.

review numbers

automatic
chapter num-
bers

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 contains our performance analysis of ZIL-LWB with a PMEM SLOG, making the case for a PMEM-specific ZIL implementation. The design and implementation of ZIL-PMEM is then presented in Chapter 4, followed by its evaluation in Chapter 5. We conclude with a summary of our findings in Chapter 6.

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-writecache [] is a new Linux Device Mapper target that synchronously persists writes to PMEM and performs asynchronous writes-back to the origin block device in the background. Dm-writecache is relevant to ZIL-PMEM because both systems expose a PMEM-accelerated virtual block device. For this use case, a write-back cache can be more space efficient because it only needs to store the most recent version of a block. In contrast, ZIL-PMEM logs each block write (smallest block size for ZVOLs is 4k) separately and does not coalesce or delta-encode log entries. ZIL-PMEM protects data integrity through checksums whereas dm-writecache fully relies on hardware error correction and detection, reported via `memcpy_mcsafe()`. Once ZVOL modifications are txg-synced to the main pool, they benefit from ZFS’s strong data redundancy mechanisms such as *raidz*. With dm-writecache, the origin block device driver must handle data redundancy if desired, e.g., though another Device Mapper target such as *dm-*

backref

ref block de-
vice compari-
son

verify claim

raid. We compare the performance of the two implementations in Section X.Y): dm-writecache delivers on the promise of low latency but shows significant lock contention with multiple threads, failing to saturate PMEM bandwidth and performing significantly worse than ZIL-PMEM.

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.

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

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.

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

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

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 inte-

grating 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

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. `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 pseudo-code in Figure 2.2 illustrates the construction of the commit list by `zil_commit`. Figure 2.3 provides an example for a single transaction group.

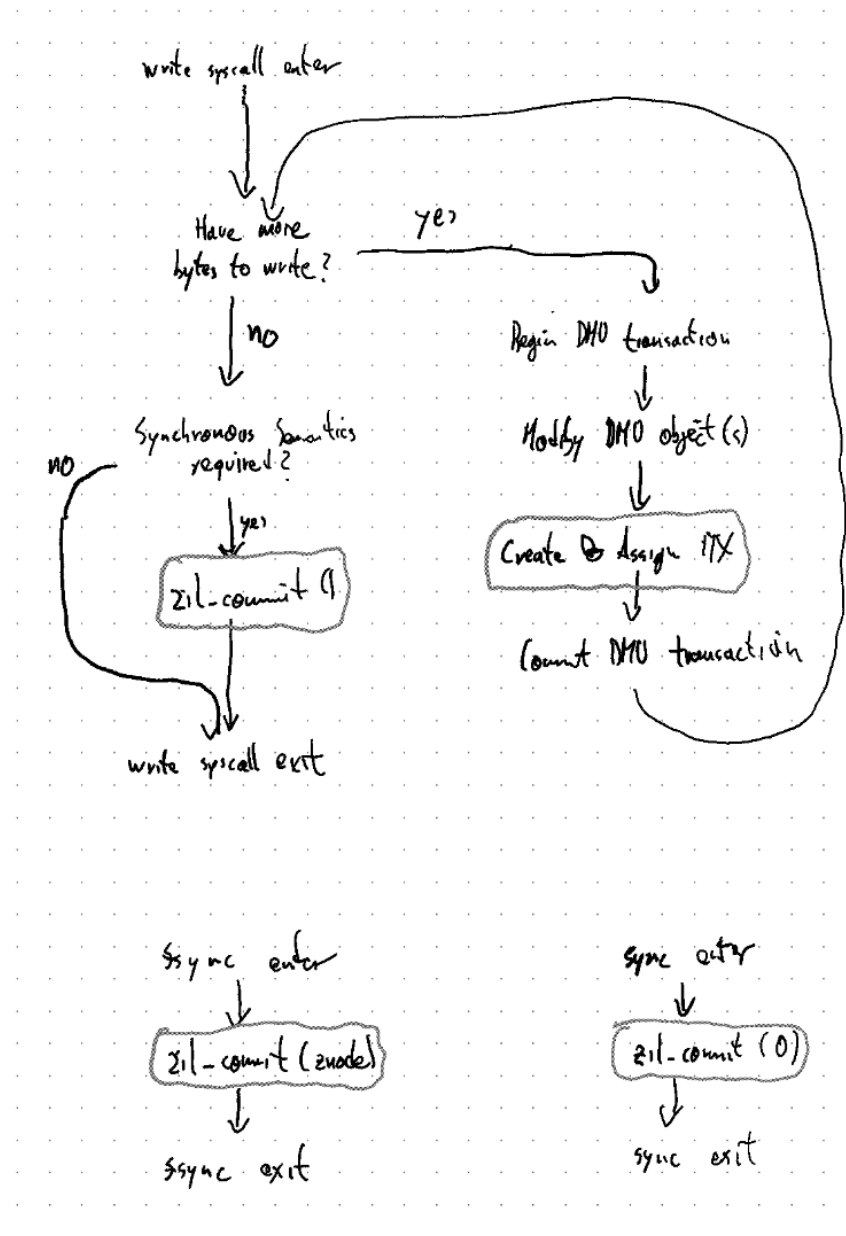


Figure 2.1: ZIL API usage on the write path in `write()`, `fsync()`, and `sync()` system calls.

```

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: Pseudo-code that illustrates the work done by `zil_commit()`.

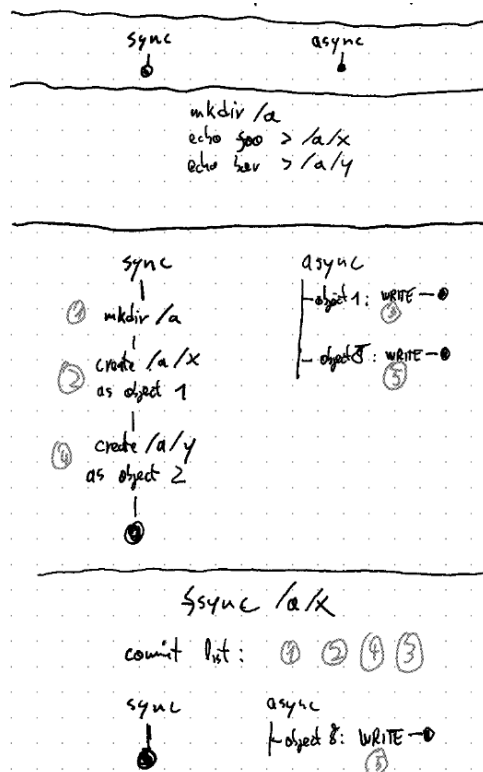


Figure 2.3: An example of how an `itxg` is filled with ITXs and how `zil_commit` drains it into a commit list. Note that this example only covers the case where all ITXs were assigned for the same `txg`.

Chapter 3

Why ZIL-LWB Is Slow On PMEM

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. _____

blindtext

Chapter 4

Design & Implementation

In this chapter we present the design and implementation of ZIL-PMEM in OpenZFS. We start with an outline of the project goals and limitations, followed by a high-level overview of the design. Each of the subsequent sections then presents one orthogonal concept in more detail, including implementation aspects where appropriate. Readers should be familiar with the technical background on PMEM in Linux and OpenZFS given in Section .

ref

4.1 Project Scope

4.1.1 Requirements

The following features are success criteria for the ZIL-PMEM design.

Coexistence ZIL-PMEM must coexist with ZIL-LWB due to limited availability of PMEM and 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 of log entries.
- 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 replay must be able to resume. Resumed replay must continue to respect logical dependencies of log entries.

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

Low Latency The latency overhead of ZIL-PMEM compared to raw PMEM device latency 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, ZIL-PMEM should achieve the following speedups for threads that perform synchronous I/O operations on separate CPU cores:

1 private dataset per thread Always near-linear speedup.

1 shared dataset

ZPL filesystem No speedup.

ZVOL Up to linear speedup, depending on workload.

Maximum Performance On Intel Optane DC Persistent Memory Whereas battery-backed NVDIMMs have existed for decades, Intel Optane DC Persistent Memory is currently the only broadly available flash-based persistent memory product. We develop ZIL-PMEM on this platform and want to determine the maximum performance that can be achieved with it.

CPU-Efficient Handling Of PMEM Bandwidth Limits The PMEM programming model dictates that I/O wait time is spent on-CPU — typically at a memory barrier instruction or because the CPU has exhausted its store or load buffer capacity. 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

reference

proof

review terminology; need proof?

at ten or more CPU cores. In contrast, DRAM shows a near-linear increase in bandwidth to 60 GB/s at 15 threads. [Yan+20, fig.4] In practice, these results imply that on-CPU time is wasted as soon as the system writes to PMEM at higher than maximum device bandwidth. Since ZIL-PMEM shares PMEM among all datasets in the pool, we expect that such overload situations 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 presented in this chapter 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 accomodate them without major changes.

Support For OpenZFS Native Encryption The ZIL-PMEM design presented in this section does not support 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 our design can be trivially extended with similar mechanisms.

4.1.3 Limitations

The following features were deliberately left out of our design. More experimentation and experience with ZIL-PMEM will benecessary 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]. Given the latency

precise ref

contribution of ZIL-PMEM to overall syscall latency (ref 8), variations of this magnitude would be significant.

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

Only Works With SLOGs ZIL-PMEM only works with PMEM SLOGs, not for a zpool with 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 logs large write records' data 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 is a 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 before we need to worry about hardware capacity.

move to background?

ref ix systems
truenas M

4.2 Overview

This section summarizes the ZIL-PMEM design at a high level. The subsequent sections then provide more detail on each abstractions' design and implementation.

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. ZIL-LWB is the current ZIL's persistence implementation. It uses the SPA's metaslab allocator to allocate log-write blocks (LWBs) from the storage pool with a bias towards SLOG devices. ZIL-PMEM disables metaslab allocation for its PMEM SLOG device and uses the PMEM space directly.

The PMEM space is partitioned into fixed-size segments. Each segment has a corresponding in-DRAM data structure called *chunk* that tracks the segment's kernel-virtual base address and length. Chunks are organized in a pool-wide in-DRAM data structure called *PRB*. PRB implements a high-performance, scalable, crash-consistent, data-corruption-checking and garbage-collected write-ahead log. It stores log entries in the PMEM chunk's segments as a contiguous append-only sequence. Full chunks are stashed away and become reusable for logging when the youngest entry's transaction group has synced to the main pool. After a system crash, a new instance of PRB is instantiated with the same chunks (segments) that the pre-crash PRB instance used. The new instance scans each chunk for log entries that need to be replayed. This process is called *claiming*. The entries' physical location in PMEM is irrelevant for replay. Instead, each entry stores sufficient metadata to determine whether an entry needs to be replayed, to detect missing entries, and to find a deterministic replay order. Once claiming is complete, the new PRB instance is ready for replay and logging by datasets. It only uses those chunks for logging that do not have a replay claim.

PRB is a pool-wide data structure but the ZIL is written and replayed on a per-dataset basis. For this purpose, PRB defines a data structure called HDL that holds the per-dataset PRB state, such as the log's GUID, dependency-tracking state and replay position. All interaction that is scoped to a dataset happens through the HDL, not PRB. Whereas PRB/HDL implements PMEM persistence, it is the HDL user's responsibility to persist HDL state in the main pool.

The keystone to the ZIL-PMEM architecture is the per-dataset in-DRAM struct `zilog_t`. Before the introduction of ZIL kinds, `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 data structures that track which log entries need to be persisted when sync 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. The subclass for ZIL-PMEM

is `zilog_pmem_t`. It is a thin wrapper around HDL’s logging and replay methods. `zilog_pmem_t` persists HDL state in the dataset’s ZIL header. The `zil_pmem` module — which defines `zilog_pmem_t` — is the component that integrates PRB/HDL into ZFS: it allocates the chunks, constructs the PRB and creates and destroys HDLs in synchrony with their datasets.

ZIL-PMEM Architecture

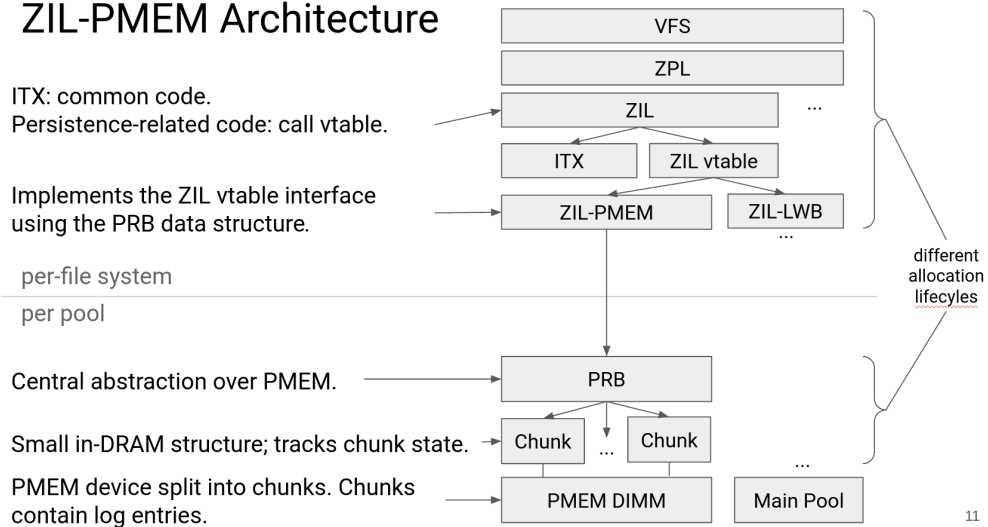


Figure 4.1: Overview of the system architecture as described in this section.

actual figure

4.3 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 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 refactored 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.

4.3.1 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 type in that union. We maintain the same size as the original `zil_header_t` by using one 64 bit word of padding, thereby simplifying the on-disk format migration. Figure 4.2 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 4.2: The ZIL header structs (in DRAM and on disk) before and after the introduction of ZIL kinds.

4.3.2 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 of ZIL records that need to be persisted before `zil_commit` can return, and in what order these log records must be replayed.
- 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 majority of functions can 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 (see Section 4.3.1 and Figure 4.2). Figures 4.3 and 4.4 illustrate the changes described in this section.

4.3.3 Changing ZIL Kinds

In the design overview (sec. ??) we introduce ZIL kinds as a pool-scoped variable. The pool's ZIL kind implicitly according to 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.

The rule is re-evaluated whenever a SLOG is added or removed, i.e., during pool creation and on `zpool add` and `zpool remove`. The idea behind this approach is to keep administration simple — a defining paradigm for ZFS: ZIL-PMEM should be automatically activated if hardware and software support it and not be used otherwise.

To accomplish this behavior, the following steps must happen in the same transaction group as the SLOG addition or removal that triggered the change of ZIL kind:

1. Stop using the ZIL API and wait until all active calls to it have finished.
2. Wait for all ZIL entries to become obsolete by waiting for the open txg to sync.

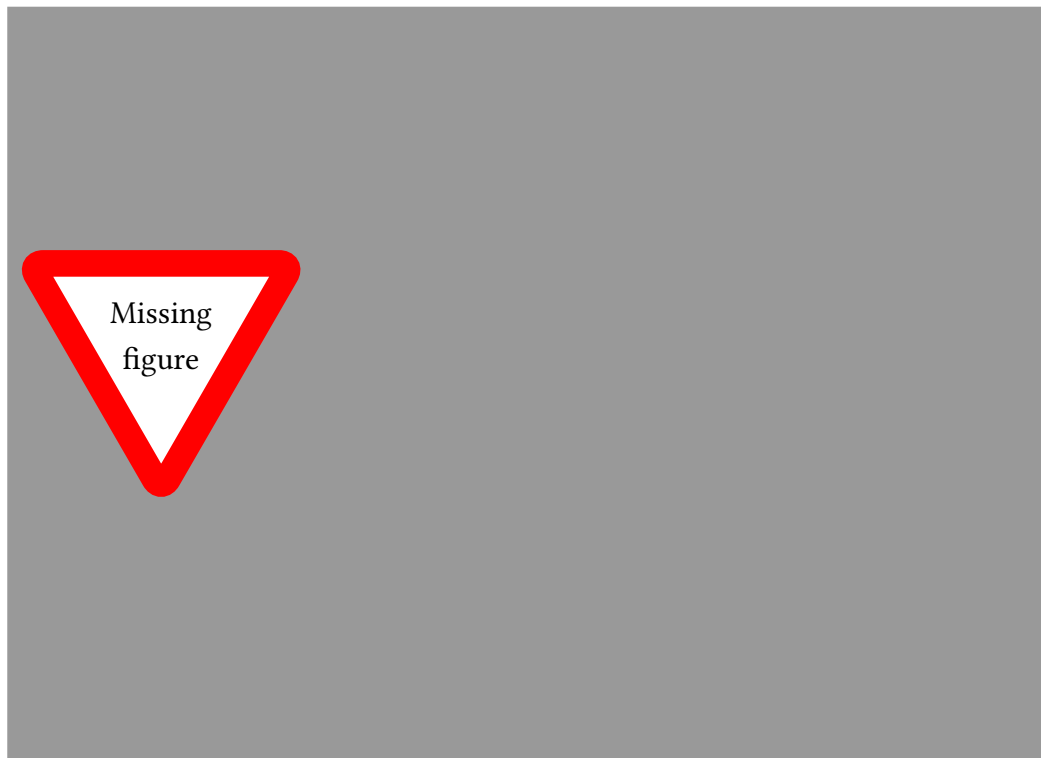


Figure 4.3: `zilog_t` before and after the introduction of ZIL kinds, visualized as a class diagram.

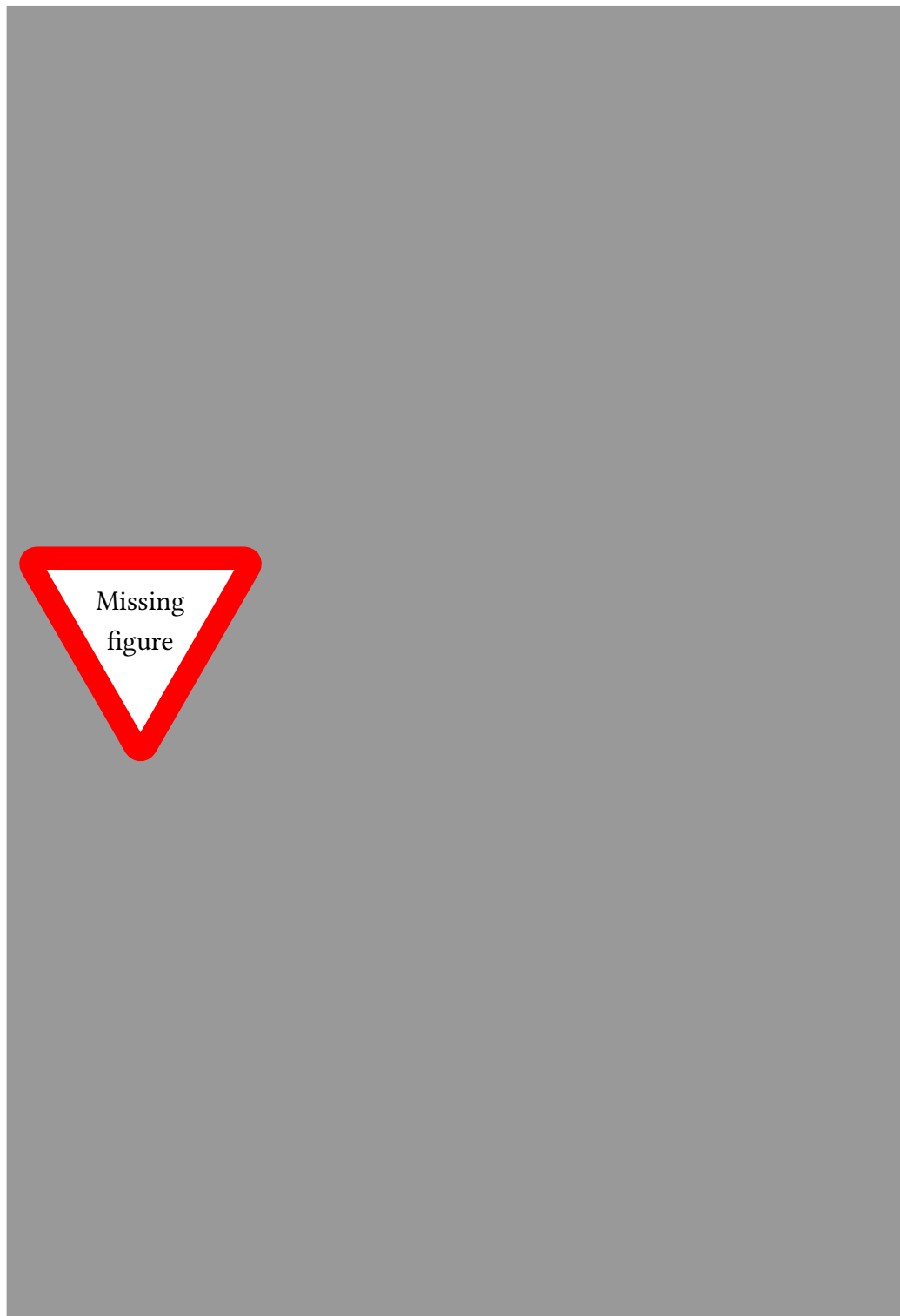


Figure 4.4: The refactoring steps described in Section 4.3.2 by example of the `zil_itx_assign()` and `zil_commit()` APIs.

3. `zil_close()` and subsequently `free()` the currently open `zilog_t` instances.
4. Replace all datasets' ZIL header with the new ZIL kind's default header. `zh_kind` is now set to the new ZIL kind.
5. Allocate the new `zilog_t` instances and `zil_open()` them. The allocation routine uses `zh_kind` to select the (new) `vtable`, allocates `zlv_t_alloc_size` bytes for the new `zilog_KIND_t` and runs the ZIL kind specific constructor.

Due to time constraints we have **not yet implemented** this procedure. As a stop-gap solution, we add 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. There is 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.

4.3.4 ZIL-LWB Suspend & Resume

The ZIL API provides the `zil_suspend` and `zil_resume` functions. `zil_suspend` pauses all ZIL activity and ensures that all log entries are obsolete. `zil_resume` reverts the state to normal operation. Compatibility code for versions of ZFS prior to the *fast snapshots* feature must use them for snapshots. For newer pool versions, the only consumer `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. Once the ZIL resumes and starts allocating LWBs, the metaslab allocator transparently uses other SLOGs or the main pool devices.

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 4.3.3) 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.

4.3.5 ZIL Traversal & ZDB

Whereas ZIL replay is 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 implementation details to its callers such as the concept of LWBs and blockpointers. 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 analysis of `zil_parse`'s 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 traverse data structures updated 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 does not actually use the ZIL data because the data integrity checks are implemented transparently in the ZIO read pipeline. 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 implement the following change as a precursor to the refactoring of `zilog_t` described in Section 4.3.2. 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. Finally, we change the `dmu_traverse` API so that callers must be explicit about ZIL-LWB traversal which we believe will ease maintainability in the future. We `zil_destroy_sync` in `old_synchronous_dataset_destroy`, leaving `spa_load_verify` as the only consumer of the `dmu_traverse` API that needs to traverse ZIL-LWB chain.

4.3.6 ZIL-LWB-Specific Callbacks

There are several callbacks from modules outside of ZIL-LWB. Since these callbacks are made through static function calls — not function pointers — they are part of the ZIL's public API and need to be considered for our ZIL kind refactoring. However, none of them are necessary for ZIL-PMEM. Therefore, we prefix the callback functions with `zillwb_`, move them to `zil_lwb.c` and ignore them in the remainder of this thesis. For maintainability, future work should replace

these statically dispatched callbacks with dynamic callbacks through function pointers to increase decoupling. We list the APIs in question so that this issue can be addressed in future work:

zil_lwb_add_txg Necessary to keep the in-DRAM representation of an LWB alive when writing WR_INDIRECT blocks. **Note:** we can only afford to ignore this function because ZIL-PMEM does not support 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.

4.3.7 Considered 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 moving the dispatch into ZIL kind specific code one API layer up so that `zilog_t` and ITX would be completely specific to ZIL-LWB. The API layer that we considered 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 ?? provides an example of how a write system call uses this API. 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

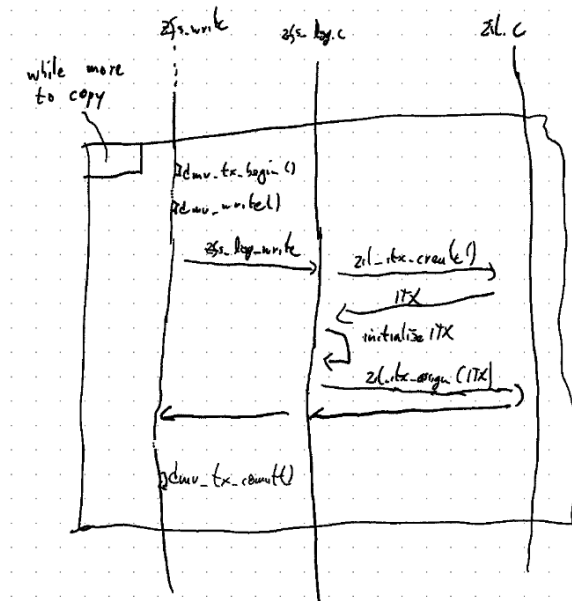


Figure 4.5: Sequence diagram of the APIs involved in creating the ITXs for a write system call.

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 hardware (see Section 4.3.3). 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.

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

ref

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 ??).

- ref
3. The PRB's feature to log entries for the same dataset in parallel would not have been possible if ZIL-PMEM would be constrained to LWBs. Our ITXG bypass for ZVOLs (see Section ??) shows how we can use this PRB feature to increase scalability without compromising on crash consistency guarantees.
- check,ref

4.3.8 Summary

We believe that our design for ZIL kinds introduces ZIL kind specific behavior at the layer that allows for sufficient flexibility in the implementation without adding unnecessary abstractions and maintenance burden. The interface defined by the vtable is a clean abstraction although some design questions (see 4.3.3 and 4.3.4) as well as some LWB-specific APIs (see 4.3.5 and 4.3.6) 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.

"sondern"?

4.4 PMEM-aware SPA & VDEV layer

In this section we describe how we add explicit support for PMEM to ZFS. Our primary goal is to disable the SPA's metaslab allocation on PMEM SLOG devices so that ZIL-PMEM can use the allocatable space as PRB chunks.

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 DAX APIs. Its value 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 log 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.

To disable metaslab allocation on PMEM SLOGs, we introduce a new allocation class called `exempt`. When adding a log device to the pool that is `is_dax`, the `zpool` command assigns it to the `exempt` class instead of the `log` class. This alone is sufficient to prevent allocation from PMEM SLOGs because the `exempt` allocation class by convention is never used. The allocatable space on the device is thus never read or written through ZIO and available for use by ZIL-PMEM. Note that reads and writes to VDEV labels are located outside of the allocatable space. Updates to VDEV labels continue to be made using block device IO (`zio_read_phys` and `zio_write_phys`).

need to explain this? short analogy: partition header

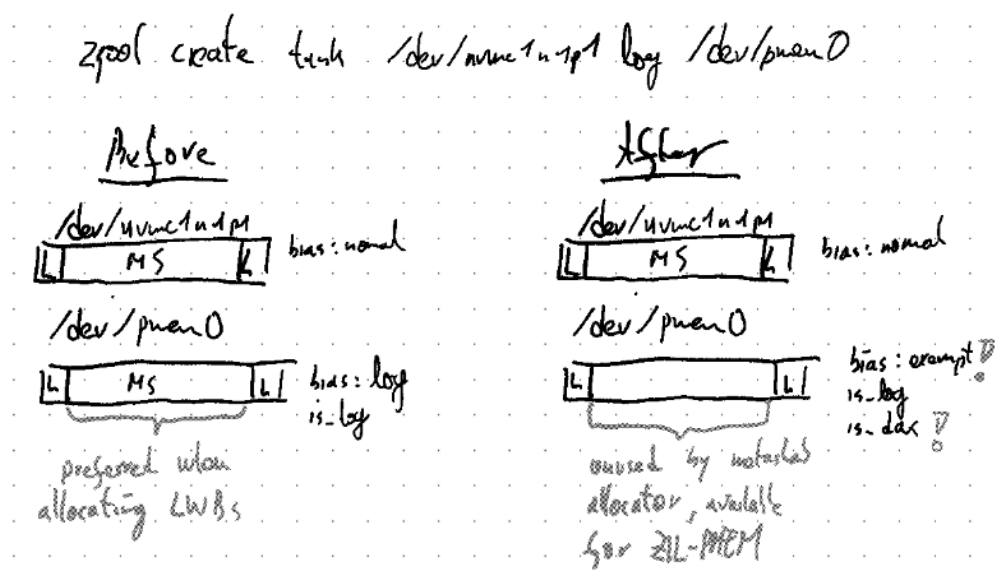


Figure 4.6: Comparison of the zpool layout before and after the addition of the `is_dax` vdev attribute and the `exempt` allocation class for PMEM SLOGs.

Chapter 5

PRB/HDL

In this section we describe the PRB data structure which implements the bulk of ZIL-PMEM’s functionality. PRB abstracts shared physical PMEM into virtual logs (HDLs) 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 approach. Sections 5.4, 5.5, 5.6, 5.7, and 5.8 progressively refine the mental model of an individual virtual log and how it is replayed. Subsequently, we describe the storage substrate that maps the per-dataset virtual logs to the shared PMEM space. Sections 5.9 and 5.10 present the data structures that we use to manage the PMEM space and store log entries. In Section 5.11 we describe how these data structures are traversed during log recovery. Section 5.12 then presents the low-latency and CPU efficient design of the write path. Finally, Section 5.13 provides an overview of the PRB API that is consumed by *zil_pmem.c*.

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 $T_{i_{txg}}$. 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 purpose of the ZIL is to bridge 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. If the system call has synchronous semantics, the system call handler invokes `zil_commit` before returning to userspace. `zil_commit` uses the `itxg` data structure to determine the sequence of log records that needs to be written out to persistent storage. This sequence is called the *commit list*. At this point, the ZIL-kind specific implementation takes over: it is responsible for concatenating the current `zil_commit` call's commit list to previously written commit lists. The logical structure of the log, regardless of ZIL kind, is thus a sequential chain of log records that is extended by `zil_commit`.

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 younger than *precrash-txg* (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 given by the commit list, but entries for transactions before or in *precrash-txg* ($T_{itxg} \leq \text{precrash-txg}$) 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.

Note that the DMU transactions during replay only re-enact the logical changes but do not imitate the original dataset. 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:

Pre-crash-txg The *pre-crash-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.

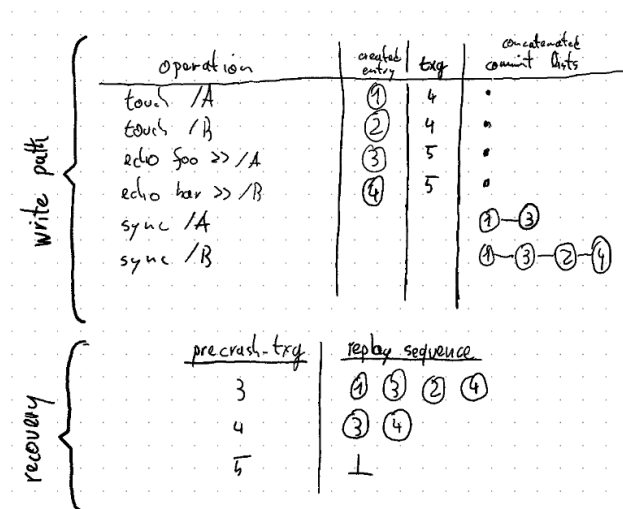


Figure 5.1: Example for a commit list with entries from different files and txgs. We show the replay sequences for different pre-crash-txgs. The replay callback must only be invoked for the entries that log changes in txgs younger than the pre-crash-txg.

5.2 Analysis

We identify the following requirements for the persistence layer of any ZIL kind:

- During normal operation, it acts as a write-only sink for log records.
- A log record is the encoded representation of the change that was made in a single DMU transaction. The representation is shared between all ZIL kinds. Log entries are always scoped to a single dataset.
- 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.

towards?

We put the following requirements on the **storage substrate** that stores the log entries:

- On the write path, the overhead that it adds to the raw PMEM latency 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. *HDL* builds on top of *PRB* a virtual log that puts individual entries into a mostly sequential structure. After a system crash, the *HDL* provides a replay facility that recovers the replayable entries of a dataset, 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.

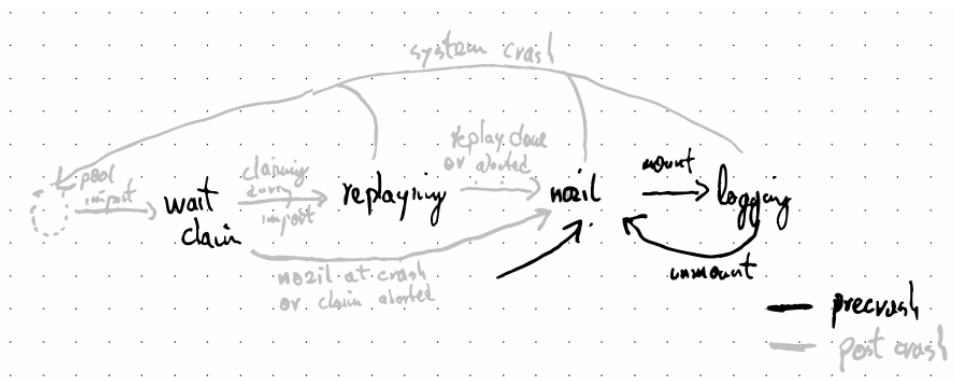


Figure 5.2: TODO

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.

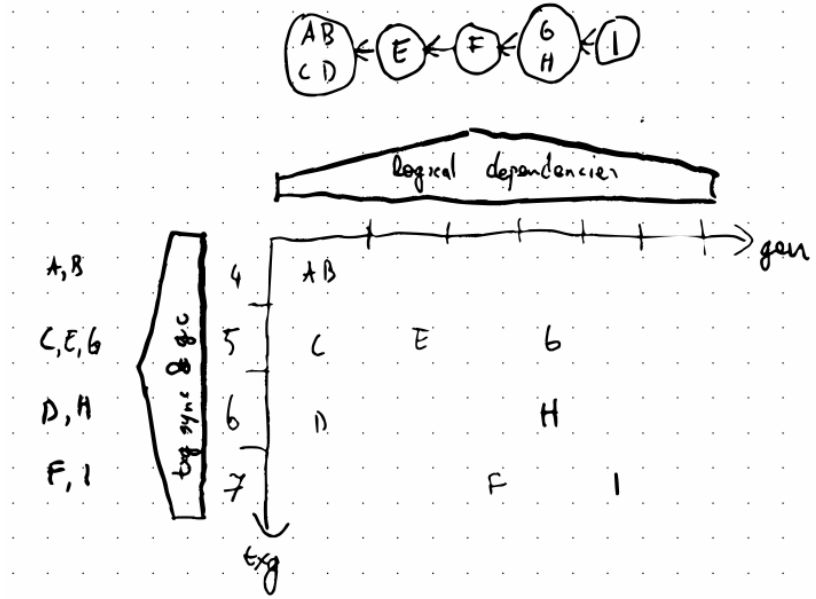


Figure 5.3: Structure of a single HDL's log as described in the previous paragraph.

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, it scans the PRB for entries that belong to the HDL and determines a *replay sequence* S .

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

The *replay order* on S defines in which 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*.

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 such a missing entry E_m with $E_m.\text{txg} > \text{precrash_txg}$ is missing, any entry E_d that logically depends on E_m ($E_m < E_d$) and is for an unsynced txg ($E_d.\text{txg} > \text{precrash_txg}$) 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 ($E_p.\text{txg} > \text{precrash_txg}$) 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 $E.C_{ctx_i}$

$$E.C_{ctx_i} := \#\{D : D.gen < E.gen \wedge D.ctx = ctx_i\}$$

stores the number of entries that were written to the log since its inception, for a DMU transaction with $ctx\ ctx_i$, 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.

The per- ctx 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 $ctx\ t_{syncd}$ has been synced, PRB garbage-collects all entries \mathcal{S}_{gc} that are obsolete:

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

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 $E.C_{ctx_i}$ in the entry metadata whose transaction groups ctx_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 ctx s and b) E only depends on entries D with $D.gen < E.gen$. Since there are only three possible unsynced ctx s at any time (*open*, *quiescing*, *syncing*), we only need to store three (ctx_i, C_{ctx_i}) tuples per log entry. In fact, since ctx s never skip a number, we only store the value of ctx_{open} and recompute

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

We use the same algorithm to compute the counters on both 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 $R_i := (ctx, ctr), i \in \{0, 1, 2, 3\}$. When an entry E is written or visited, it modifies the counter in the row with index $I_T := T \bmod 4$ where $T := E.ctx$. We distinguish the following conditions:

compact encoding is just an idea, not in impl yet

language, help

- 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 $T < R_{I_T}$, 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*.

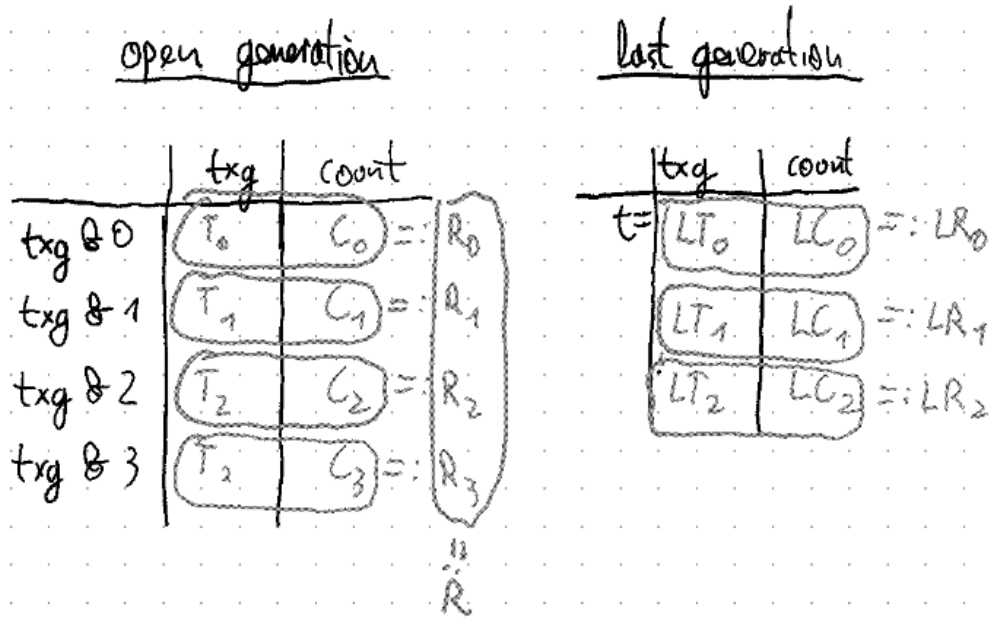
Regardless of whether a new generation was started or not, we always update $E_{last} \leftarrow E$ and increment the counter in the *live table* as described in the previous paragraph after evaluation the rules above.

Since the counters $E.C_{txg_i}$ only count up to but not including $E.gen$, they same for all entries in generation $E.gen$. 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 $R_i \in \text{last table}$.
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.

figures need to be adjusted to new formulation of algorithm



Rules to compute LR_0, LR_1, LR_2

Let

$$R_x.T := T_x ; \quad LR_x.T := LT_x ;$$

$$LR_x = R_x \Leftrightarrow LT_x = R_x \wedge LC_x = C_x$$

in

$$LR_0 := R_i \text{ where } T_i \text{ is } \max(T_0, T_1, T_2, T_3).$$

$$LR_i := \begin{cases} (0, 0) & \text{if } LR_{i-1}.T \leq 1 \\ R_j & \text{if } \exists R_j \in R : R_j.T = LR_{i-1}.T - 1 \\ (0, 0) & \text{otherwise} \end{cases}$$

Figure 5.5: Visualization of the *live* and *last* table structures, and a graphical illustration of the algorithm that is used to determine the counters that are stored in the entries of a generation. Note that we also provide an Example in Figure 5.7.

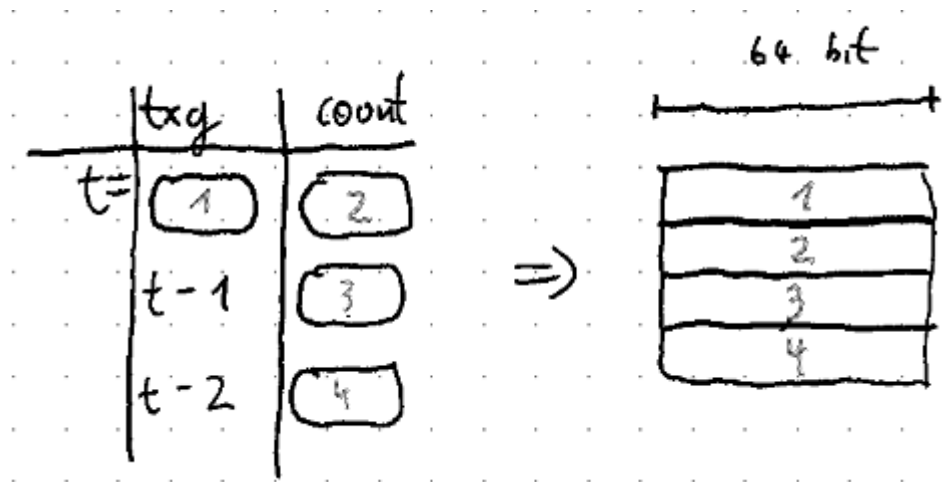
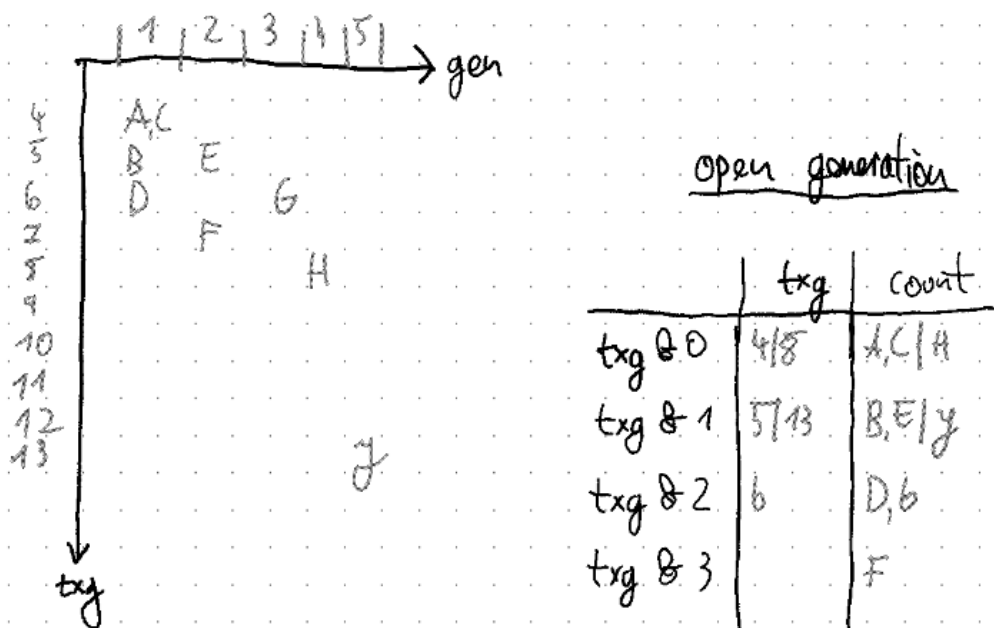


Figure 5.6: Space-efficient encoding of the counters.



Emitted Counter Tables by generation:

generation	(txg, counter) for		
	txg	txg-1	txg-2
1	0 0	0 0	0 0
2	6 1	5 1	4 2
3	7 1	6 1	5 2
4	7 1	6 2	5 2
5	8 1	7 1	6 2

Each row represents one generation's counter table

	txg	count
t	6	1
t-1		1
t-2		2

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 some times but stay in the corridor of at most three a corridor of 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 "|" 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 that produce
- 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.

the impl currently doesn’t check for collisions, need to mention that?

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*:

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_{last} The last entry that was replayed.

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.
 - 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.

unless it's at the end of the chain, but we don't need to mention that here...

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.

have a section on DRAM usage? better have it be one big section where we discuss this...

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 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 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 ??.

Chunks can only be modified by one thread at a time. To support parallel writes, 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. Chunks that are not associated with a commit slot are always on exactly one of the following lists:

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 writers that need a new chunk for their 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 recount*. Once claiming is done and *txg sync* triggers the first garbage collection cycle, any chunk on this list that has a zero recount is emptied and moved to the *free list*. The list is also scanned for zero recounts 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's has transitioned to state *nozil* on disk and crashed inbetween, the garbage-collected entries might be missing when resuming replay.

Figure 5.8 visualizes the above chunk ownership transitions. Figure 5.9 provides a corresponding example.

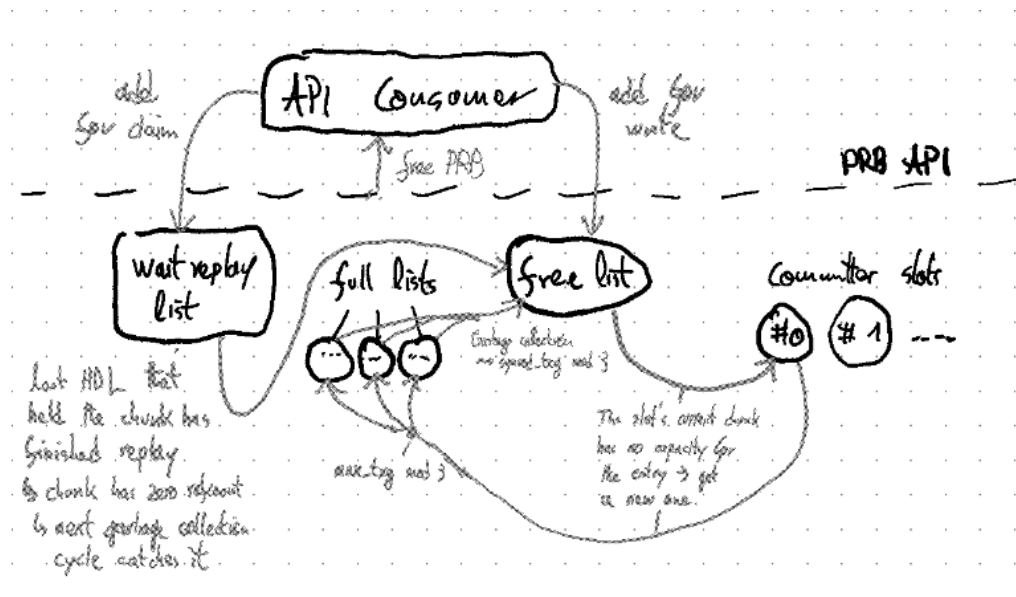


Figure 5.8: The different owners of a chunk and the events that cause ownership transitions.

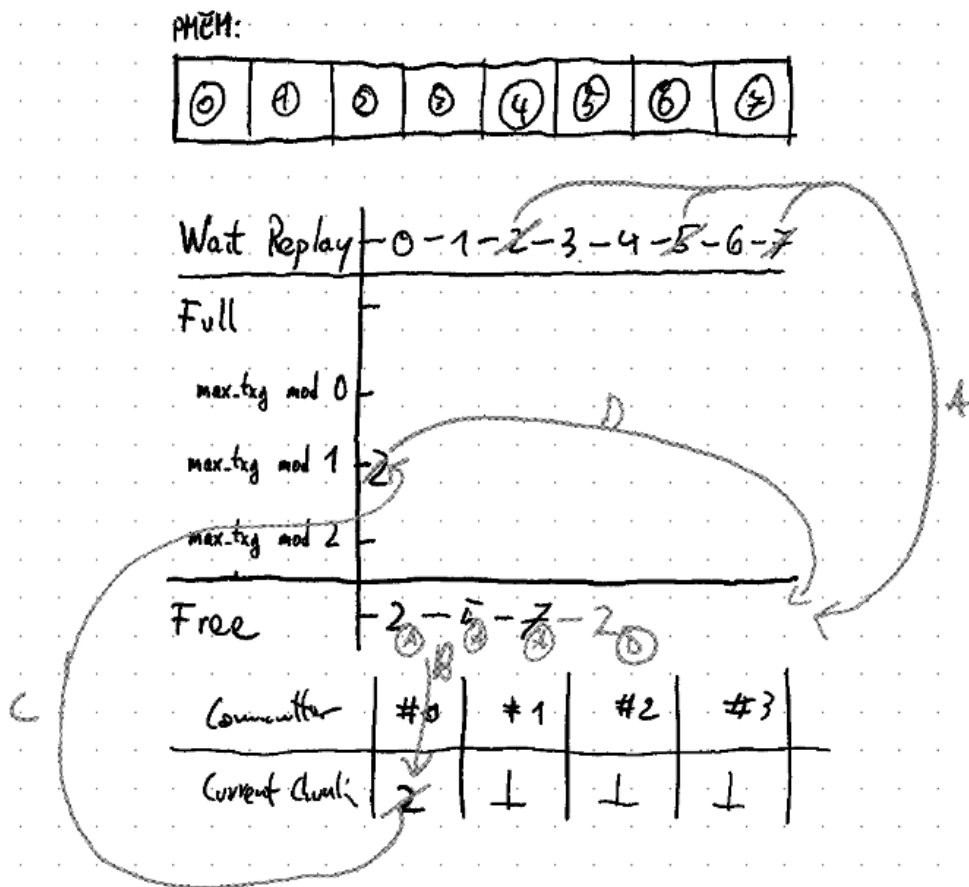


Figure 5.9: Example for chunk ownership transitions. The PRB is constructed with all eight chunks added as `_for_claim`. Claiming (A) determines that chunks 0, 1, 3, 4, and 6 need to remain on the *wait replay* list because they contain entries for HDL logs that need replay. Chunks 2, 5, and 7 only contain 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.

parametrize
checksum algo-
rithm, store in
HDL?

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 5.12.

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.

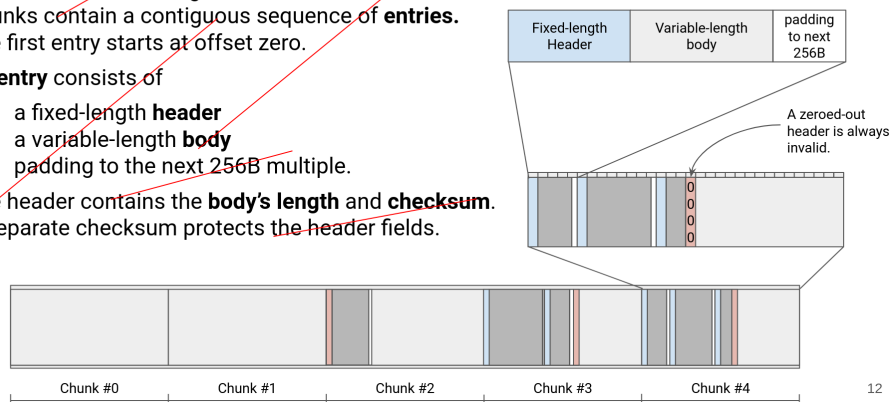


Figure 5.10: Example chunk layout. Note that although we partition the PMEM space very regularly in this example, the PRB consumer is free to use variable-sized chunks.

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. For example, a zero *log GUID* is defines as invalid.
- 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 different cases:

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 do 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.

not the current
api, but it's
shorter...

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 5.12.3. 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 ch_len >= sizeof(entry_header_t)
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));

    validate_header_checksum(eh);

    if eh.log_guid == 0 {
        // invalid header terminates the sequence
        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: Write Path

A thread that writes to a HDL needs to perform the following tasks.

- Find a target chunk using the *commit slot* mechanism. (Section 5.12.1)
- 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 7).

linux percpu
counter pre-
cise?

Commit Slot Acquisition & 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 — typically at a memory barrier instruction or because the CPU has exhausted its store or load buffer capacity. This is problematic from a CPU utilization perspective. For example, 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.

review terminology; need proof?

5.12.1 Commit Slots

Commit slots are our abstraction to enable multiple threads to write entries concurrently. The goal is to grant up to $n_{\text{committers}}$ 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 *acquires a commit slot* $S \in 0, 1, \dots, n_{\text{committers}} - 1$. Each thread that is simultaneously committing gets a different commit slot. If no commit slot is available, the function to acquire the commit slot blocks. Associated with each commit slot is a PMEM chunk which we refer to as *open chunk*. A thread that acquires a commit slot is allowed 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 acquires 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 acquisition and release procedures.

i think that's a common term?

Procedure For Commit Slot Acquisition:

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 7), 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:

Blast Radius Entries are packed into a small number of chunks, even more so if we checked all *open chunk* to implement some form of *best fit* placement of entries. Packing is beneficial for space efficiency because the *spread* between mininum and maximum txg of the entries in a chunk is small, enabling timely garbage collection. But 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.

Cache Efficiency Our aquisition procedure deterministcally picks the lowest availble commit slot. It is thus very likely that chunks bounce around CPU

chapter?

maybe we should change the design to always use non-blocking?

cache or CPU affinity?

cores, without taking cache affinity into consideration. Note that, due to the use of non-temporal store instructions when writing to PMEM, this only impacts the chunk DRAM object.

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

5.12.2 HDL-Scoped Metadata

check matches

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.

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 7), 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

check content

Section 7.

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.

5.12.3 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 ch.pos 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.

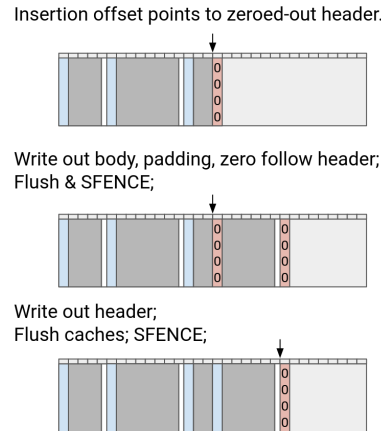


Figure 5.11: The crash-consistent append operation to a PMEM chunk. The traversal algorithm reaches the existing entries at all times and reaches the new entry only after its body and header have been written completely.

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 $\text{sizeof}(lr_write_t) + \text{data} + \text{padding} = 4096 + 192 + 64 = 4352$ bytes large. Assuming 2 GiB/s write bandwidth for a single Optane DIMM ([Yan+20]), the write time for the entry body is ~ 2 us. In contrast, the derived latency for a 256 byte write at that rate is ~ 0.12 us. If we use this value as an approximation for the cost of the `sfence`, its latency contribution is only $\frac{0.12}{0.12+2} = 5.6\%$. Note that we have not systematically evaluated the effect.
- 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.
-

need to prove this? don't have time...

should we do this? if so, need also evaluate on different pmem configurations

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

5.13 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.

5.13.1 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_teardown_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_teardown_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_tgx` is the pool's first new transaction group. For HDLs in state *logging*, `pool_first_tgx - 1` becomes the *precrash_tgx*.

`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`.

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

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

5.13.2 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 unplayable 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 dis-
cards the unre-
playable part of
the log

5.13.3 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 5.12.2. 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 parallel to the same HDL. Note that writers for different HDLs need not coordinate at all.

contrary?

language?

Garbage 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

ZIL-PMEM

In this section we describe how we use PRB/HDL to implement ZIL-PMEM as a new ZIL kind.

Chapter 7

ITXG Bypass For ZVOL

In this section we describe a performance optimization for ZIL-PMEM that allows for parallel log writes to a single ZVOL.

```
typedef struct zilpmem_replay_state_phys {
    uint64_t claim_txg;
    eh_dep_t resume_state_active;
    eh_dep_t resume_state_last;
} zilpmem_replay_state_phys_t;

typedef struct eh_dep {
    uint64_t eh_last_gen;
    prb_deptrack_count_pair_t eh_last_gen_counts[TXG_CONCURRENT_STATES];
} eh_dep_t;

typedef struct prb_deptrack_txg_seq_pair {
    uint64_t dtp_txg; /* 0 ==> invalid pair */
    uint64_t dtp_count; /* 0 ==> invalid pair */
} prb_deptrack_count_pair_t;
```


Chapter 8

Evaluation

8.1 Correctness

We evaluate the correctness of our work through unit and integration tests.

8.1.1 PRB

8.1.2 ZIL-PMEM

8.2 Performance

We evaluate the performance of ZIL-PMEM with a variety I/O benchmarks.

8.2.1 Write Performance

4k Random Sync Writes

Application Benchmarks

8.2.2 Replay Performance

Chapter 9

Related Work

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).
- [] *Linux 5.0 Compat: SIMD Compatibility · Openzfs/Zfs@e5db313*. GitHub. URL: [/openzfs/zfs/commit/e5db31349484e5e859c7a942eb15b98d68ce5b4d](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).
- [] *Ndctl-Inject-Error Man Page - Ndctl - General Commands*. URL: <https://www.mankier.com/1/ndctl-inject-error> (visited on 04/13/2021).
- [] *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 Sankaranarayanan 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.
- [SDM413] Intel Corporation. *Intel 64 and IA-32 Architectures Software Developer’s Manual*. Vol. 3. 325462-074US. Apr. 2021. Chap. 4.1.3 Paging-Mode Modifiers.
- [Vol+14] Haris Volos et al. “Aerie: Flexible File-System Interfaces to Storage-Class Memory.” In: *Proceedings of the Ninth European Conference on Computer Systems*. 2014, pp. 1–14.
- [XS16] Jian Xu and Steven Swanson. “ $\{\$NOVA\}$: A Log-Structured File System for Hybrid Volatile/Non-Volatile Main Memories.” In: *14th $\{\$USENIX\}$ Conference on File and Storage Technologies ($\{\$FAST\}$ 16)*. 2016, pp. 323–338.
- [Xu+17] Jian Xu et al. “NOVA-Fortis: A Fault-Tolerant Non-Volatile Main Memory File System.” In: *Proceedings of the 26th Symposium on Operating Systems Principles*. 2017, pp. 478–496.
- [Yan+06] Junfeng Yang et al. “Using Model Checking to Find Serious File System Errors.” In: *ACM Transactions on Computer Systems (TOCS)* 24.4 (2006), pp. 393–423.

- [Yan+20] Jian Yang et al. “An Empirical Guide to the Behavior and Use of Scalable Persistent Memory.” In: 18th {USENIX Conference on File and Storage Technologies ({FAST 20). 2020, pp. 169–182. ISBN: 978-1-939133-12-0. URL: <https://www.usenix.org/conference/fast20/presentation/yang> (visited on 04/26/2021).
- [YCH19] Takeshi Yoshimura, Tatsuhiko Chiba, and Hiroshi Horii. “EvFS: User-Level, Event-Driven File System for Non-Volatile Memory.” In: 11th {USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 19). 2019. URL: <https://www.usenix.org/conference/hotstorage19/presentation/yoshimura> (visited on 10/02/2020).
- [Zha+21] Wenhui Zhang et al. “ChameleonDB: A Key-Value Store for Optane Persistent Memory.” In: (2021).
- [ZHS19] Shengan Zheng, Morteza Hoseinzadeh, and Steven Swanson. “Zigurat: A Tiered File System for Non-Volatile Main Memories and Disks.” In: *17th {USENIX} Conference on File and Storage Technologies ({FAST} 19)*. 2019, pp. 207–219.