

September 23-26, 2019 Santa Clara, CA

#### Volatile Use of Persistent Memory

Usha Upadhyayula Intel Corporation



## Agenda

- Motivation
  - Why use persistent memory as volatile memory
- Exposing persistent memory as volatile memory
  - File-Backed memory
  - As a NUMA node
- Call to action

2



## MOTIVATION

## Intel® Optane<sup>™</sup> DC Persistent Memory Operating Modes

#### Memory Mode (volatile)

- Data Placement
  - Application does not have control over data placement
    - DRAM and Intel® Optane<sup>™</sup> DC Persistent Memory
- Ease of adoption: No Code Changes
- Performance slower than DRAM

#### APPLICATION VOLATILE MEMORY POOL DRAM AS CACHE OPTANE PERSISTENT MEMORY Rights Reserved

#### App Direct Mode (persistent)

- Data Placement
  - Control over data placement
- Ease of Adoption:
  - Need Code Changes
- Native hardware latencies



## Motivation

- Applications need large memory capacity
- Don't need persistence
- Need control over data placement
  - DRAM and other storage tiers
- Native latencies of persistent memory



#### **SD**©

# File-Backed Memory: Using Memory Mapped Files

### **File-Backed Memory**

- Memory mapped files
  - Requires DAX aware file system
    - XFS, EXT4, NTFS
  - Bypasses file system page cache
- Fastest IO path possible
  - No Kernel code or Interrupts
- Code changes required for Load/Store Access



SDI(9

## **Memory Allocation Challenges**

- malloc/free don't work on memory mapped files
- Improve ease of use
  - stdlib like API to allocate memory

#### **Solution: libmemkind**







2019 Storage Developer Conference. © Insert Your Company Name. All Rights Reserved.

## File-backedTemporary file created

& memory mapped on a persistent memoryaware file system

Supports malloc/free

Allocations not persistent

PMEM KIND

interface

- Temporary file deleted when the application exits
- Need simple modifications to the applications

#### libmemkind- How it Works



#### SD @

## **Memkind API**

#### KIND Creation

- Fixed & variable size heap
- Automatic KIND detection
  - Static and dynamic KINDs
- Destroy KIND
- KIND HEAP Management
  - Allocate
  - Free
  - Usable size
  - Detect KIND
- KIND Configuration Management
  - Usage Policy, Set Path, Set Size

#### S D @

#### Code Walkthrough

## Create & Allocate from Different KINDs

```
int main(int argc, char *argv[])
```

```
struct memkind *pmem_kind = NULL;
memkind_create_pmem("/mnt/pmem", PMEM_MAX_SIZE, &pmem_kind);
//allocate in DRAM
char * ptr_default = (char *)memkind_malloc(MEMKIND_DEFAULT, size);
//allocate in file backed "Kind" of memory
char * ptr_pmem = (char *)memkind_malloc(pmem_kind, size);
//Free allocated memory
memkind_free(MEMKIND_DEFAULT, ptr_default);
memkind_free(pmem_kind, ptr_pmem);
memkind destroy kind(pmem kind);
return 0;
```

https://github.com/memkind/memkind/tree/master/examples/pmem\_and\_default\_kind.c

2019 Storage Developer Conference. © Intel® Corporation. All Rights Reserved.

SDO

## **Create with Defined and Unlimited Size**

- // Create first PMEM partition with specific size
- struct memkind \*pmem\_kind = NULL
- err = memkind\_create\_pmem(/mnt/pmem, PMEM\_MAX\_SIZE, &pmem\_kind);
- // Create second PMEM partition with unlimited size
- err = memkind\_create\_pmem(/mnt/pmem, 0, &pmem\_kind\_unlimited);
- // Destroy both PMEM partitions
- err = memkind\_destroy\_kind(pmem\_kind);
- err = memkind\_destroy\_kind(pmem\_kind\_unlimited);

https://github.com/memkind/memkind/tree/master/examples/pmem\_kinds.c

SD(®

## **Using KIND Configuration API**

// Create new configuration

struct memkind\_config \*pmem\_cfg = memkind\_config\_new(); memkind\_config\_set\_path(pmem\_cfg, "/mnt/pmem/"); memkind\_config\_set\_size(pmem\_cfg, 1024 \* 1024 \* 64); memkind\_config\_set\_memory\_usage\_policy(pmem\_cfg, MEMKIND\_MEM\_USAGE\_POLICY\_CONSERVATIVE); // Create pmem\_kind with configuration set memkind\_create\_pmem\_with\_config(pmem\_cfg, &pmem\_kind); memkind\_config\_delete(pmem\_cfg); memkind destroy kind(pmem kind);

https://github.com/memkind/memkind/tree/master/examples/pmem\_config.c

2019 Storage Developer Conference. © Intel® Corporation. All Rights Reserved.

15

#### **Code Samples**

File Name	Description	
pmem_kinds.c	Create and destroy PMEM kind with defined or unlimited size.	
pmem_malloc.c	Allocate memory and the possibility to exceed PMEM kind size.	
pmem_malloc_unlimited.c	Allocate memory with unlimited kind size.	
pmem_usable_size.c	View the difference between the expected and the actual allocation size.	
pmem_alignment.c	Use memkind alignment and how it affects allocations.	
pmem_multithreads.c	Use multithreading with independent PMEM kinds.	
pmem_multithreads_oneki nd.c	Use multithreading with one main PMEM kind.	
pmem_and_default_kind.c	Allocate in standard memory and file-backed memory (PMEM kind).	
pmem_detect_kind.c	Distinguish allocation from different kinds using the detect kind function.	
pmem_config.c	Use custom configuration to create PMEM kind.	
pmem_free_with_unknown	Allocate in standard memory & file-backed memory (PMEM kind),	
<u>kind.c</u>	and free memory without needing to remember which kind it belongs to.	
https://github.com/memkind/memkind/tree/master/examples.		

#### **Libmemkind - Recap**

- General purpose library
- Uses jemalloc for memory management
- Provides stdlib like interface
- Uses file-backed memory
  - Memory map temp files

#### An Efficient & Scalable Cache for Volatile use of Persistent Memory

#### Cache Design Challenges: Fragmentation & Scalability

- Challenges
  - Random memory allocation sizes
  - Long runtime durations
  - Large memory capacities
- Problem
  - Failure to allocate a contiguous chunk of memory although the requested chunk is availble
- Existing solutions
  - Compacting GC (Java, .NET)
  - Defragmentation (Redis, Apache Ignite)
  - Slab allocation (memcached)





2019 Storage Developer Conference. © Intel® Corporation. All Rights Reserved.

SD(®

## Iibvmemcache – Caching Solution for Large Memory Capacity

- Main Features
  - Custom memory allocator
    - Extent-based
      - Reduces fragmentation
    - Control over allocations
    - Improves space efficiency
  - Buffered LRU
    - Delivers scalability
      - Works with Large Memory Capacities
  - Uses memory mapped files





**SD**<sup>®</sup>

## **Extent Based Allocation**

- Similar to file system extents
  - Extent = contiguous set of blocks
- Multiple non-contiguous blocks allocated
- Non-contiguous allocations appear as a single allocation





#### SD©

21

#### **Cache Design - Buffered LRU**

- Added a wait-free ring-buffer
  - buffers the list-move operations
- List only needs to get locked only during eviction or when the ring-buffer is full.



SD©

#### Design Aspects: libmemkind vs libvmemcache

	libmemkind (PMEM)	libvmemcache
Allocation	Dynamic allocator	Extent based
Scheme		
Purpose	General purpose	Key-value store optimized
		for large memory
		capacities
Fragmentation	Apps with random size	Minimized
	allocations/deallocations that	
	run for a longer period	

#### Cache Design - Lightweight, embeddable, in-memory caching

SD®

VMEMcache \*cache = vmemcache\_new("/mnt/pmem", VMEMCACHE\_MIN\_POOL, VMEMCACHE\_MIN\_EXTENT, VMEMCACHE\_REPLACEMENT\_LRU);

vmemcache delete(cache);

#### Libvmemcache - Recap

#### libvmemcache

- get/put APIs
- optional replacement policy
- configurable extent size
- Works with terabyte-sized in-memory workloads
- High space utilization



#### Intel® Optane<sup>™</sup> DC Persistent Memory Module

https://github.com/pmem/vmemcache

# Persistent Memory as an Extension of DRAM

#### **Persistent Memory as DRAM extension**

- New feature in Linux\* kernel 5.1
  - Dev\_DAX\_KMEM Config Option
  - Binds persistent memory device to kernel
  - Appears as a separate NUMA node
- libmemkind support
  - memkind\_malloc with new static KIND
    - MEMKIND\_DAX\_KMEM



SD @

## **Call To Action**

- Get started with persistent memory programming
  - <u>https://github.com/pmem/pmdk</u>
  - Libvmemcache
  - <u>https://github.com/memkind/memkind</u>
- Learn more about volatile usages
  - Accelerate your Apache Spark with Intel® Optane<sup>™</sup> DC Persistent Memory
  - Pmem-redis
- Join the development efforts
  - pmem.io <u>http://pmem.io/</u>
  - Intel Persistent Memory Developer Zone
- Send us your feedback

#### Let's Innovate As a Community

2019 Storage Developer Conference.  $\textcircled{\sc 0}$  Intel  $\textcircled{\sc 0}$  Corporation. All Rights Reserved.





2019 Storage Developer Conference. © Insert Your Company Name. All Rights Reserved.

## Cache Design - Scalable replacement policy

- Performance of libvmemcache was bottlenecked by naïve implementation of LRU based on a doubly-linked list.
- With 100st of threads, most of the time of any request was spent waiting on a list lock...
- Locking per-node doesn't solve the problem...



#### Intel<sup>®</sup> Optane<sup>™</sup> DC Persistent Memory **Operating Modes**





#### **AFFORDABLE & LARGE** VOLATILE MEMORY CAPACITY

PERSISTENT MEMORY