

## 2021 OFA Virtual Workshop REMOTE PERSISTENT MEMORY ACCESS AS SIMPLE AS LOCAL MEMORY ACCESS

Tomasz Gromadzki, Software Architect (RPMem)

**Intel Corporation** 

intel

### © OpenFabrics Alliance

## WHAT IS PERSISTENT MEMORY (PMEM)?

- allow accessing data as memory
  - directly
  - byte-addressable
- the contents are non-volatile (preserved across power cycles).
- it doesn't typically replace either memory or storage
- For more information please visit:
  - pmem.io
  - Persistent Memory Development Kit (PMDK) GitHub repo

pmem.io/book





### PERSISTENT MEMORY PROGRAMMING MODEL

Storage Networking Industry Association (SNIA) standardized Software Persistent Memory (PMem) programming model

When writing data to PMem, in order for the data to be considered persistent:

• After writing to persistent memory, the software is responsible for

flushing data from CPU caches to the Persistence Domain

• Once the CPU caches are flushed of write data,

the platform guarantees the write data is persistent, should power be lost

### IMPLEMENTING THE PMEM PROGRAMMING MODEL OVER AN RDMA NETWORK (APM)



### IMPLEMENTING THE PMEM PROGRAMMING MODEL OVER AN RDMA NETWORK (GPSPM)





### RDMA accesses Intel® Optane<sup>™</sup> PMem in the same way it accesses DRAM

### **Remote PMem (RPMem)**

is about well-known technologies (like PCIe, RDMA) used in a new way

### THE NEW LIBRPMA FOCUSES ON RPMEM USABILITY



- memcpy-like API
- Hidden RDMA complexity
- Application can freely manage your PMEM all the time
- Minimum dependencies
- Enables Persistent Memory Programming Model

up to 50% RPMem source code reduction in an application that moves from libibverbs to librpma

### **LIBRPMA API**

- **Connection management** 
  - to ensure operations consistency
  - to hide RDMA complexity
- **Remote Persistent Memory Access** (RPMA)
  - Read, Write, Flush, Atomic write lacksquare
- Messaging
  - also with PMEM-backed message buffers  $\bullet$
- Memory management
  - r\_key exchange support
- Ready to incorporate RDMA Memory **Placement Extension**



### **BASIC EXAMPLE – MEMORY REGISTRATION**

# Initiator node

rpma\_mr\_reg(peer,

ptr, size,

RPMA\_MR\_USAGE\_WRITE\_SRC,

&src\_mr);



rpma\_mr\_reg(peer,

ptr, size,

RPMA\_MR\_USAGE\_WRITE\_DST

RPMA\_MR\_USAGE\_FLUSH\_TYPE\_PERSISTENT,

&dst\_mr);

### **BASIC EXAMPLE – REMOTE PERSISTENT MEMORY WRITE**

# Initiator node

rpma\_write(conn,

dst\_mr, dst\_offset,

src\_mr, src\_offset,

KILOBYTE, RPMA\_F\_COMPLETION\_ON\_ERROR, NULL);

# Target node

### **BASIC EXAMPLE – REMOTE PERSISTENT MEMORY WRITE**

# Initiator node

#### rpma\_write(conn,

dst\_mr, dst\_offset,

src\_mr, src\_offset,

KILOBYTE, RPMA\_F\_COMPLETION\_ON\_ERROR, NULL);

rpma\_flush(conn,

dst\_mr, dst\_offset,

KILOBYTE, RPMA\_FLUSH\_TYPE\_PERSISTENT, RPMA\_F\_COMPLETION\_ALWAYS, FLUSH\_ID);

https://github.com/pmem/rpma/tree/master/examples/05-flush-to-persistent

Target node

### LIBRPMA EXAMPLES COLLECTION



- > 📴 10-send-with-imm
- > 🔓 11-write-with-imm

- Connection establishment and management
- Read/write from/to DRAM/PMem
- Multiple connections (scalability)
- Atomic write
- Messaging
- Flush to persistent (both APM and GPSPM)
- Send/Write with immediate data



Server node

# Client node

rpma\_peer\_new(dev, &peer);

Common

### **CONNECTION SETUP**





# Client node



rpma\_conn\_disconnect(conn);

rpma\_conn\_next\_event(conn, &conn\_status);
assert(conn\_status == RPMA\_CONN\_CLOSED);

rpma\_conn\_delete(&conn)

rpma\_conn\_next\_event(conn, &conn\_status); if(conn\_status == RPMA\_CONN\_CLOSED) { rpma\_conn\_disconnect(conn); rpma\_conn\_delete(&conn); }

### GPSPM EXAMPLE (09-FLUSH-TO-PERSISTENT-GSPSPM)

- APM's Flush and GPSPM's Flush do the same thing but they are not semantically the same
- GPSPM Flush can be combined with other application's operations
- No common API for APM and GPSPM flushes as we want to provide messaging API for applications

### **LIBRPMA 0.9 FEATURES**

- Establishing an RDMA connection with a remote node and monitoring connection status
  - **Configurable** connection parameters
  - **Exchanging** small **information** during the connection establishment process
- Read remote memory (both persistent and volatile)
- Write to remote memory (both persistent and volatile)
- Ensuring data placement within the memory subsystem of a remote node (flush)
- Executing an atomic write (Intel platform-specific via standard aligned RDMA Write with appropriate fencing; native AtomicWrite-ready)
- Polling, in a non-blocking way, for RDMA operations completions and connection status changes
- Logging all diagnostic/status to syslog and/or on stderr
- A transparent FileSystem DAX support (on selected RNICs)
- Providing easy to use RDMA Send/Receive messaging (also with PMem-backed buffers)
- Additionally:
  - Examples collection
  - A Fio based benchmarking environment

### THE LIBRPMA LIBRARY ONE YEAR LATER

#### Use cases

- The librpma is going to be used in Ceph and SQL-like engine
- The librpma library is utilized by a number of customers as a reference solution
- The librpma library is used as a fast RPMem ramp up tool
- FIO engines for APM, GPSPM and AoF are based on librpma
  - all of these allows for easy, on-premise RPMem benchmarking
- API is stable but small tuning is still possible in response to customer feedback (0.9)
- New features under development(performance and usability improvements)
  - Separate Receive Completion Queue
  - More examples
    - Send/Recv with PMem
    - Connection errors handling

### **RPMEM BENCHMARKING ENVIRONMENT**

Fio engines (librpma APM/GPSPM, librpma AoF)

github.com/pmem/fio => github.com/axboe/fio/pull/1186

- bandwidth/latency
- remote DRAM vs RPMEM (dev-dax, fs-dax)
- numjobs, blocksize, iodepth, readwrite

#### pmem.io/rpma

- *Performance Tuning* blog describing RPMem related BIOS/OS settings
- Direct Write to PMem blog describing RPMem enabling on Intel servers
- github.com/pmem/rpma/tree/master/tools/perf
  - rpma\_fio\_bench.sh to collect performance data
    - Fio job files templates
    - remote DDIO control, local/remote NUMA control, PMem/DRAM
  - csv\_compare.py for results comparison (research, manual analysis)
  - create\_report.sh for comprehensive performance report
    - rpma\_fio\_bench.sh 192.168.1.1 all all all
    - adjustable report template

### **BENCHMARKING TOOLSET**



github.com/pmem/rpma/tree/master/tools/perf

### MORE DOCUMENTATION AVAILABLE

#### Visit

- <u>pmem.io/rpma</u> for official documentation
- github.com/pmem/rpma\_to
  - build the library
  - run the examples
  - setup the benchmarking environment

pmem.io Persistent Memory Programming					
Home About	Glossary	Documents	PMDK	ndctl	Blog
				_	
Persiste	ent Memo	ory Develop	ment K	it	
Persiste The librp	oma library	ory Develop	ment K	it	
The librp The Remote accessing pe Access (RD	ent Memo oma library Persistent M rsistent memor MA).	Pry Develop Memory Access (R ry devices on remo	ment K PMA) ( <i>librp</i> ote hosts ove	<b>it</b> oma) is a C er <b>Remote</b>	library to simplify Direct Memory
Persiste The librp The Remote accessing pe Access (RD Documentati	ent Memo pma library Persistent M rrsistent memor MA). ion is available:	emory Develop emory Access (R ry devices on remo	ment K PMA) ( <i>librp</i> ote hosts ove	<b>it</b> oma) is a C er <b>Remote</b>	library to simplify Direct Memory

You do not need neither an Intel® Optane<sup>™</sup> PMem nor RDMA-capable NIC to start examples

#### You can play with RPMA examples on any Linux desktop



2021 OFA Virtual Workshop

# **THANK YOU**

Tomasz Gromadzki, Software Architect (RPMEM)

**Intel Corporation** 



### LEGAL NOTICE AND DISCLAIMERS

This document contains information on products in the design phase of development. The information here is subject to change without notice. Do not finalize a design with this information.

No License (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software, or service activation. Performance varies depending on system configuration. Learn more at intel.com, or from the OEM or retailer.

No computer system can be absolutely secure. Intel does not assume any liability for lost or stolen data or systems or any damages resulting from such losses.

You may not use or facilitate the use of this document in connection with any infringement or other legal analysis concerning Intel products described herein. You agree to grant Intel a non-exclusive, royalty-free license to any patent claim thereafter drafted which includes subject matter disclosed herein.

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document. The products described may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.

Intel, the Intel logo, Intel Xeon, and Optane Persistent Memory are trademarks of Intel Corporation in the U.S. and/or other countries.

\*Other names and brands may be claimed as property of others.

© 2021 Intel Corporation.





### IS THE LIBRPMA ATOMIC WRITE AN ATOMIC OPERATION?

- IBTA Spec does not support atomicity for RDMA Write
- librpma implementation assumes the aligned 8-byte data won't be torn either in RDMA HW/SW stack or on the PCIe bus level
- The current implementation explicitly uses fencing to ensure ordered data processing according to the upcoming RDMA.AtomicWrite ordering rules
- The native RDMA.AtomicWrite will be used whenever and wherever it will be available