

# RPMA - Performance Report

Release CLX + MLX + CentOS 8.2 kernel 4.18

**Testing Date:** August 2021

**Performed by:**

- Łukasz Dorau (lukasz.dorau@intel.com)
- Tomasz Gromadzki (tomasz.gromadzki@intel.com)
- Jan Michalski (jan.m.michalski@intel.com)
- Oksana Sałyk (oksana.salyk@intel.com)

## Audience and Purpose

This report is intended for people who are interested in evaluating RPMA performance.

The purpose of the report is not to imply a single "correct" approach, but rather to provide a baseline of well-tested configurations and procedures that produce repeatable results. This report can also be viewed as information regarding best-known methods/practices when testing the RPMA performance.

## Disclaimer

Performance varies by use, configuration and other factors. Learn more at [www.Intel.com/PerformanceIndex](http://www.Intel.com/PerformanceIndex).

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure.

Intel technologies may require enabled hardware, software or service activation.

Your costs and results may vary.

Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

## Test Setup

### Common Configuration (both initiator and target)

Item	Description
Server Platform	S2600WFT
CPU	Intel® Xeon® Gold 6252 CPU@2.10GHz

Item	Description
Memory	12x 32GB Micron 36ASF4G72PZ-2G9E2, DDR4, 2666MT/s Total of 384GB DRAM
Persistent Memory	12x 256GB Intel® Optane™ Persistent Memory, 2666MT/s Total of 3072GB PMem 12GiB Device DAX on 6x Interleaved region
RDMA-capable NIC	100Gbps Mellanox Technologies ConnectX-5 MTU 4200 RoCEv2
Storage	Intel® SSDSCKKB48
Operating System	CentOS Linux 8 (Core)
BIOS	SE5C620.86B.02.01.0013.121520200651
Linux kernel version	4.18.0-193.28.1.el8_2.x86_64
librpma	0.9.0-452-gf5ccce3
libibverbs	1.10.30.0
libpmem	1.6.1
FIO version	<a href="https://github.com/pmem/fio/tree/rpma">https://github.com/pmem/fio/tree/rpma</a> fio-3.23-419-g79ae6
rdma-core	Version 51mlnx1 Release 1.51258
Testing Date	August 2021

## Target Configuration

Intel® Data Direct I/O Technology (DDIO) turned off on a per PCIe Root Port basis for the Mellanox NIC PCIe Root Port. This is required to allow direct writing to remote PMem via RDMA.

Item	Description
Repository	<a href="https://github.com/pmem/rpma/tree/man-pages">https://github.com/pmem/rpma/tree/man-pages</a>
Version	c907efd9d5709aa303cdaf155bc3c1f182511bf0
Applied procedure	<a href="#">Details</a>
Testing Date	August 2021

## BIOS Settings

Item	Description
Repository	<a href="https://github.com/pmem/rpma/tree/man-pages">https://github.com/pmem/rpma/tree/man-pages</a>
Version	c907efd9d5709aa303cdaf155bc3c1f182511bf0
Applied procedure	<a href="#">Details</a>
Testing Date	August 2021

Excerpt:

Item	Description
Package C-State	C0/C1 state
C1E	Disabled
CPU Power and Performance Policy	Performance
Fan Profile	Performance
Memory Configuration - Average Power Budget	18mW

Item	Description
Memory Configuration - NMV Performance Setting	Latency Optimized
Intel® Turbo Boost Technology	Enabled
Energy Efficient Turbo	Disabled
Intel® Hyper-Threading Technology	Disabled
Intel® Virtualization Technology	Disabled
Thermal Monitor	Disabled
Testing Date	August 2021

## Kernel & BIOS spectre-meltdown information

Both server systems use CentOS 4.18.0-193.28.1.el8\_2.x86\_64 kernel version available from repository with default patches for spectre-meltdown issue enabled.

BIOS on all systems was updated to post spectre-meltdown versions as well.

## Introduction to RPMA

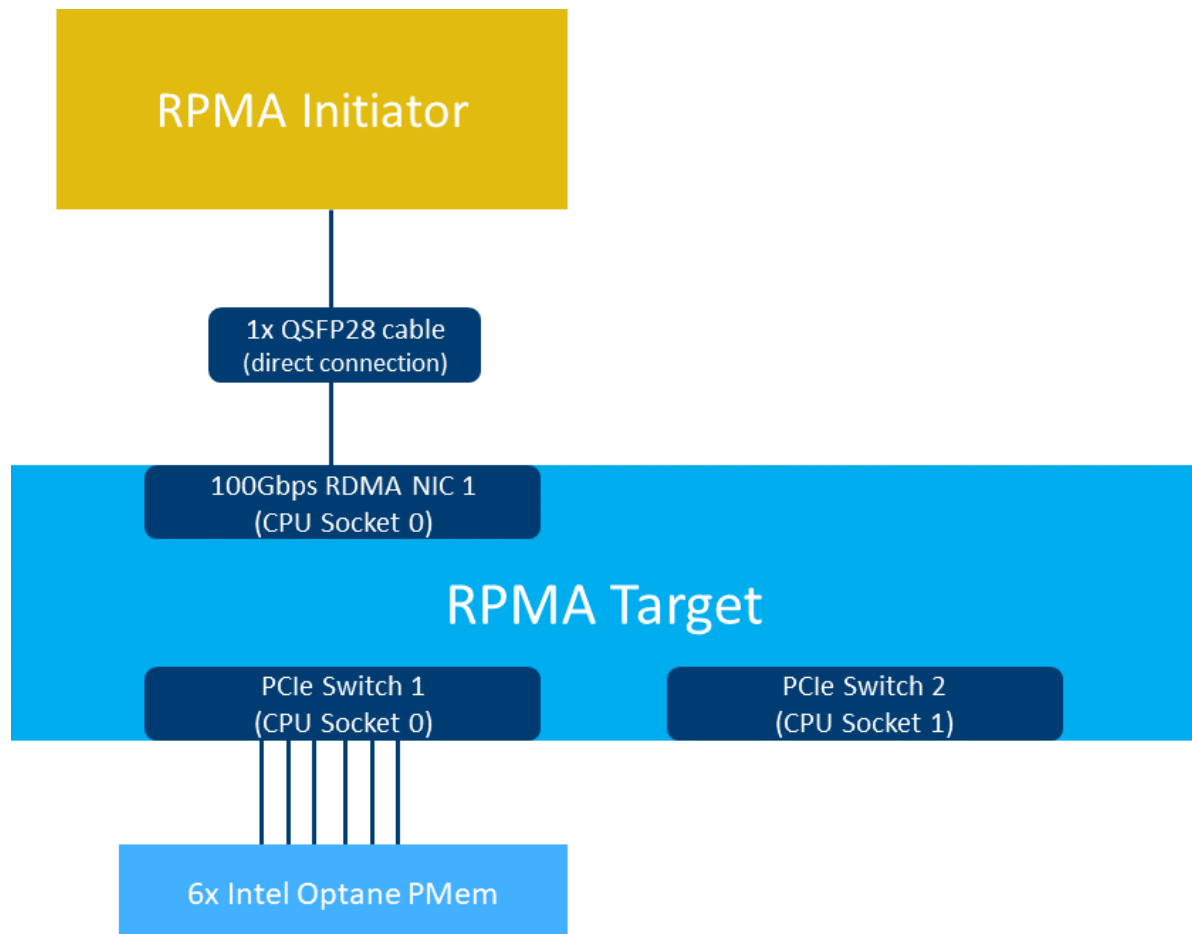
The Remote Persistent Memory Access (RPMA) library (librpma) uses Remote Direct Memory Access (RDMA) to provide easy to use interface for accessing Persistent Memory (PMem) on the remote system. It is a user-space library which may be used with all available RDMA implementations (InfiniBand™, iWARP, RoCE, RoCEv2) from various providers as long as they support the standard RDMA user-space Linux interface namely the libverbs library.

The RPMA-dedicated FIO engines are created as complementary pairs namely librpma\_apm\_client along with librpma\_apm\_server and librpma\_gpspm\_client along with librpma\_gpspm\_server. For the simplicity sake, both parts are implemented independently without any out-of-band communication or daemons allowing to prepare the target memory for I/O requests. The server engine prepares memory according to provided job file (either DRAM or PMem), registers it in the local RDMA-capable network interface (RNIC) and waits for the preconfigured (via the job file) number of incoming client's connections. The client engine establishes the preconfigured number of connections with the server. After these setup steps the client engine starts executing I/O requests against the memory exposed on the server side.

The RPMA library and any application using it (including FIO with dedicated engines) should work on all flavours of RDMA transport but it is currently tested against RoCEv2.

The FIO should be configured in a way that guarantees running all its threads and allocating all its buffers from a single NUMA node, the same the used RDMA interface is attached to, to avoid costly cross-NUMA synchronizations (e.g. using Ultra Path Interconnect).

Please see a high-level schematic of the systems used for testing in the rest of this report. The setup consists of two individual systems (one used as the initiator and one used as the target). The target system has six Intel® Optane™ Persistent Memory devices connected to the respective NUMA node. Both the initiator and the target are equipped with 100Gbps Mellanox ConnectX-5 NICs. The initiator is connected to the target system directly using QSFP28 cables without any switches. From two ports available on NICs only one is in use during the measurements.



# Read from PMem

Benchmarking the `rpma_read()` operation against various data sources (PMem, DRAM) and comparing the obtained results to standard RDMA benchmarking tools: `ib_read_lat` and `ib_read_bw`.

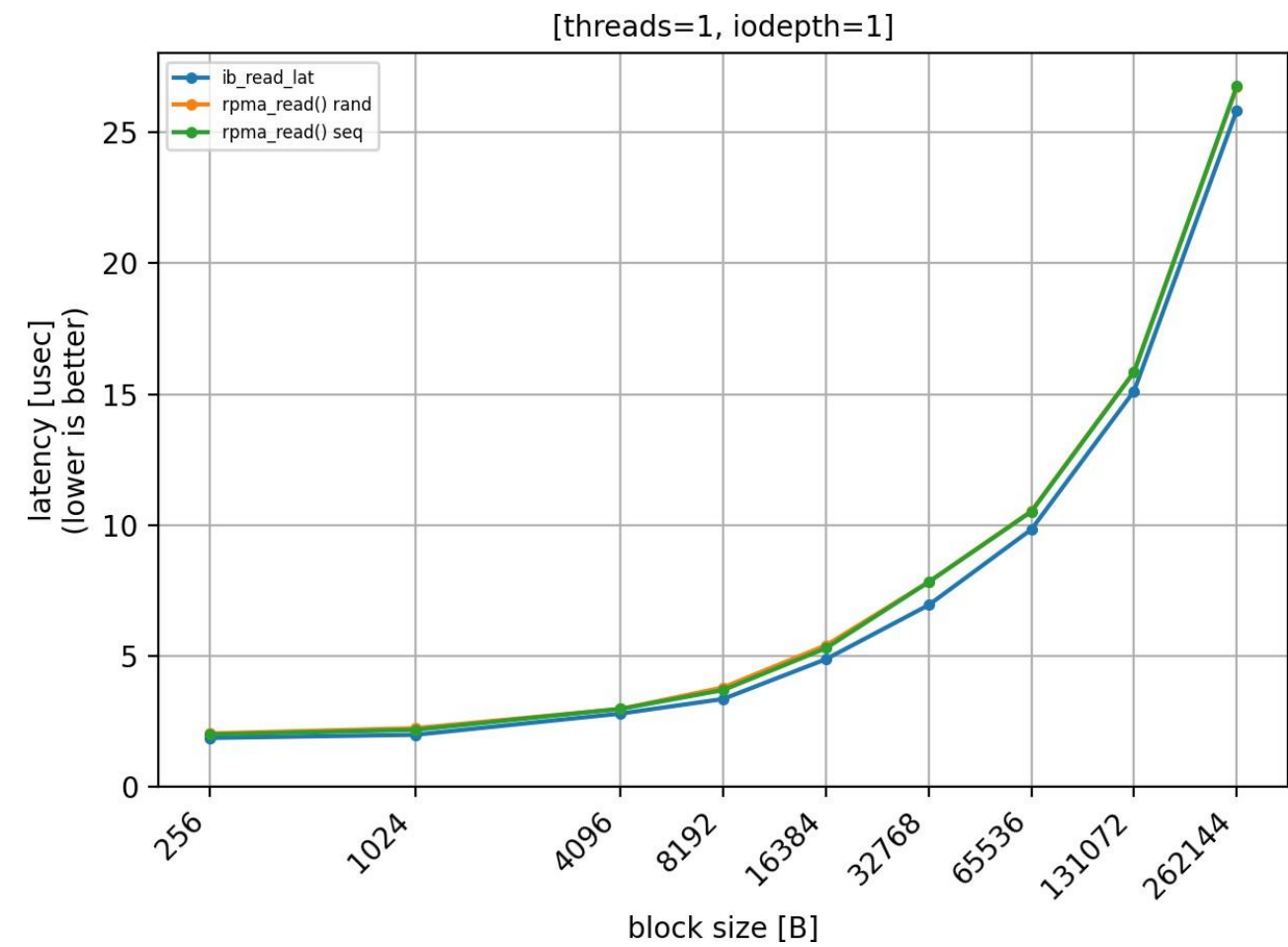
## Read from PMem: Latency

Comparing the latency of `rpma_read()` from PMem on **the RPMA Target** to the latency of `rpma_read()` from DRAM on **the RPMA Target** using the `ib_read_lat` as the baseline.

Item	Description
Server - <code>ib_read_lat</code> configuration	<code>--size \$blocksize</code>
Client - <code>ib_read_lat</code> configuration	<code>--iters \$iters --size \$blocksize --perform_warm_up \$serverip</code>
Server - FIO engine configuration	<pre>[global] ioengine=librpma_apm_server create_serialize=0 kb_base=1000 serverip=\$serverip port=7204 thread  [server] direct_write_to_pmem={0, 1} # 1 for Device DAX numjobs=1 size=100MiB filename={malloc, /dev/dax/path}</pre>

Item	Description
Client - FIO engine configuration	<pre>[global] ioengine=librpma_apm_client create_serialize=0 serverip=\$serverip port=7204 thread disable_clat=1 lat_percentiles=1 percentile_list=99.0:99.9:99.99:99.999  [client] sync=1 readwrite={read, randread} blocksize=\$blocksize ramp_time=15 time_based runtime=60</pre>

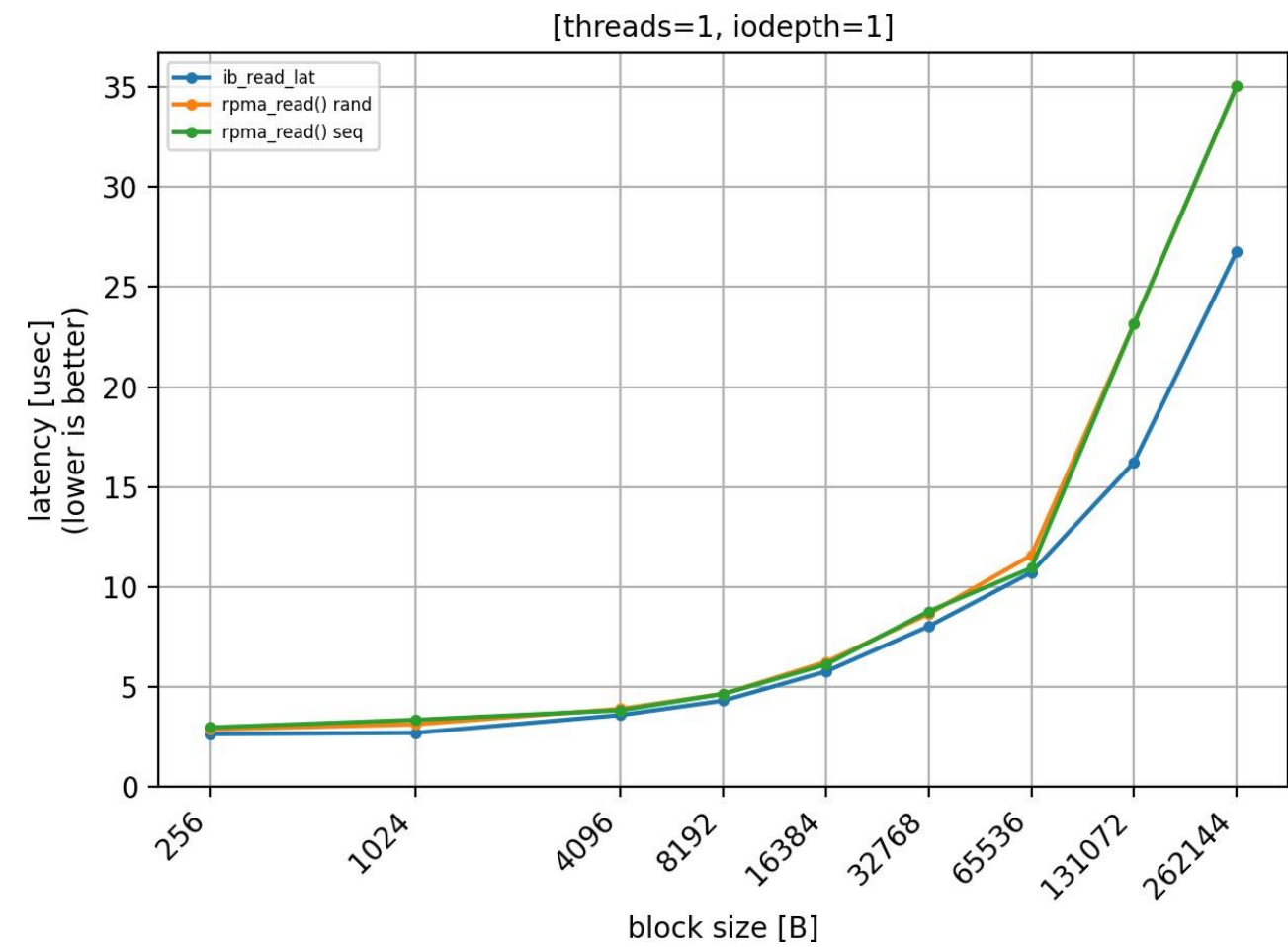
Figure 1. Latency (lat\_avg): ib\_read\_lat vs rpma\_read() from DRAM



	256	1024	4096	8192	16384	32768	65536	131072	262144
ib_read_lat	1.86	1.98	2.79	3.35	4.87	6.94	9.83	15.12	25.86
rpma_read() rand	2.03	2.23	2.97	3.79	5.40	7.82	10.53	15.86	26.75
rpma_read() seq	1.99	2.17	2.97	3.69	5.29	7.82	10.51	15.85	26.78

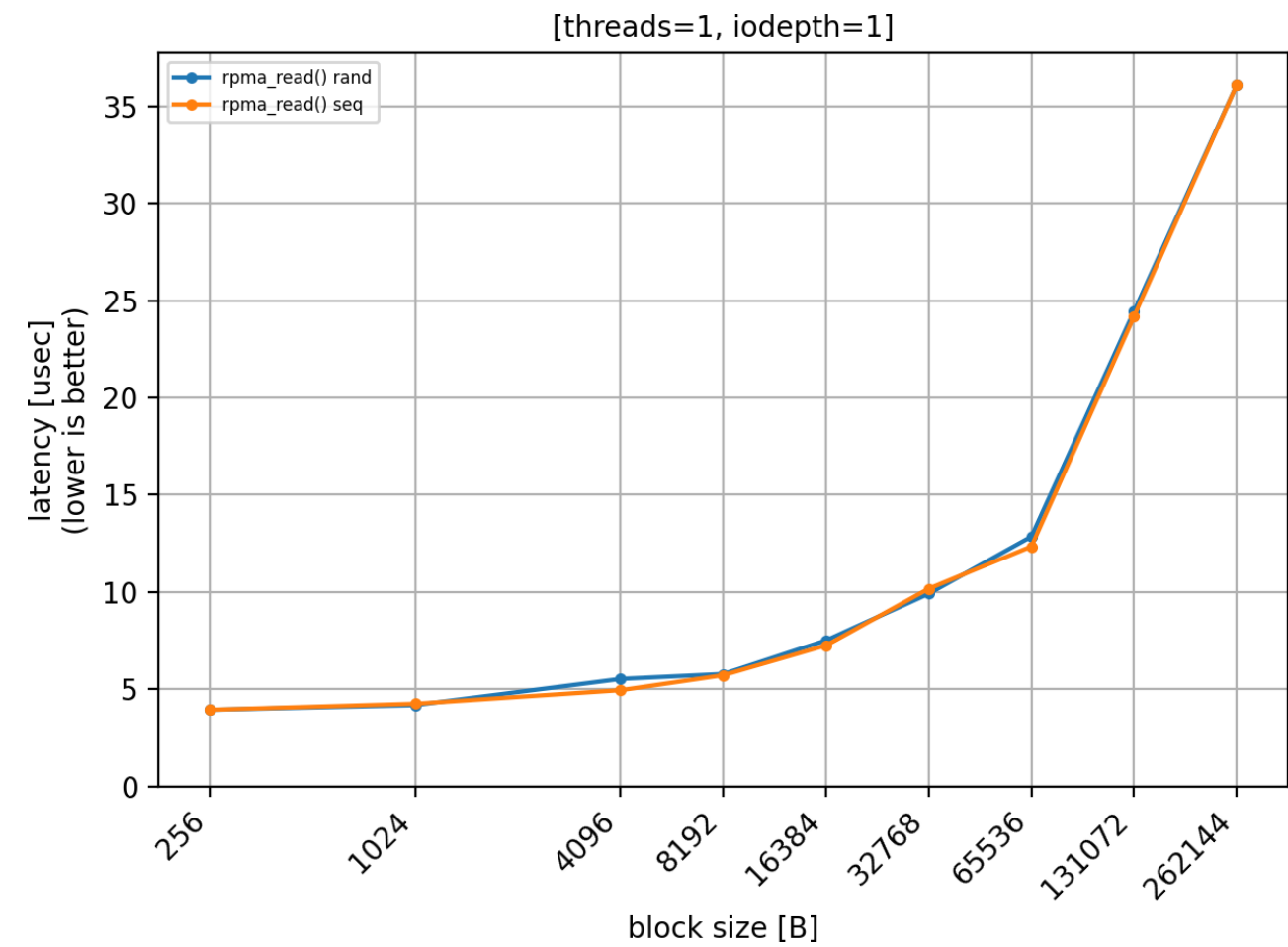


Figure 2. Latency (lat\_pctl\_99.9): ib\_read\_lat vs rpma\_read() from DRAM



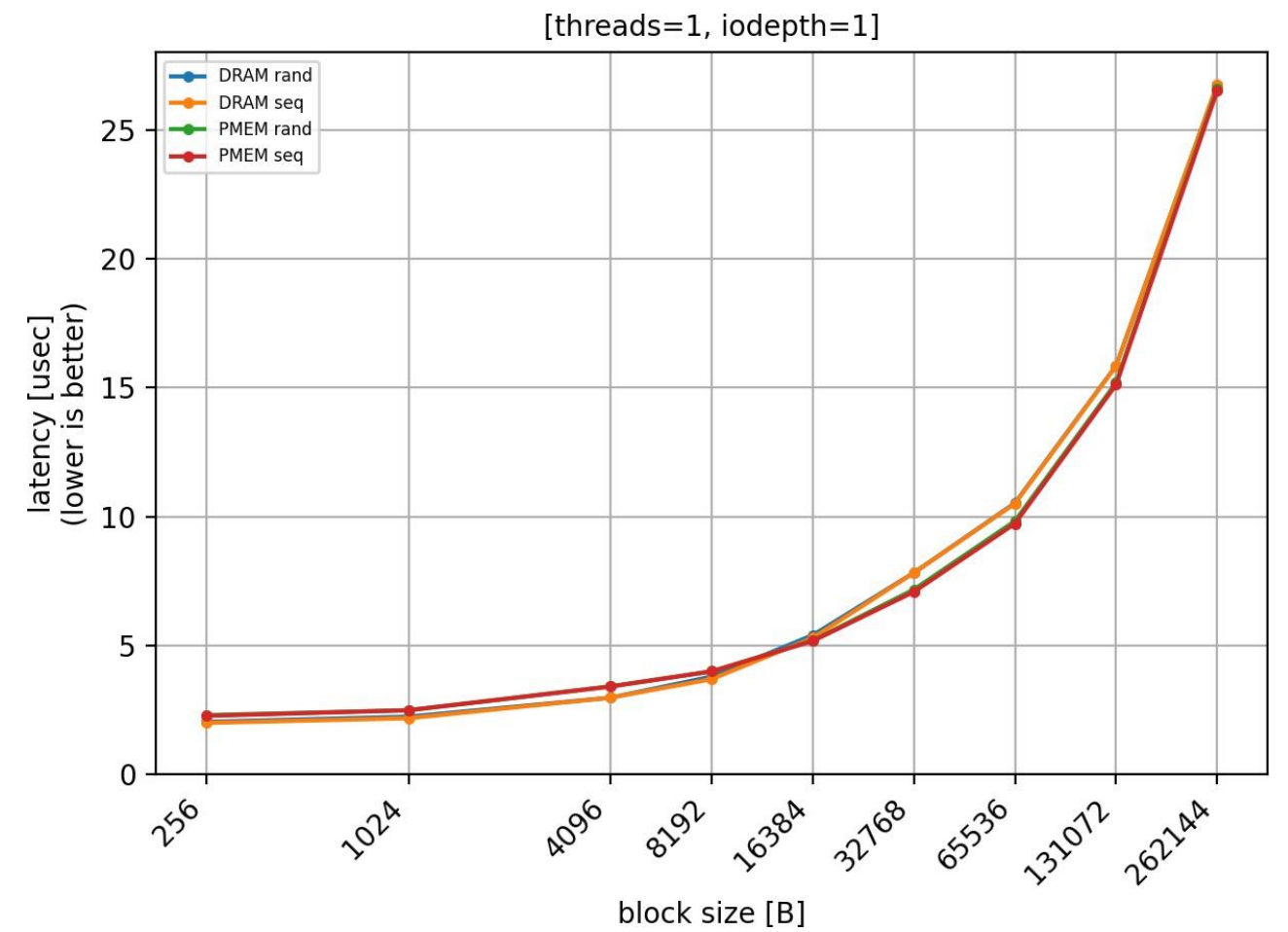
	256	1024	4096	8192	16384	32768	65536	131072	262144
ib_read_lat	2.63	2.69	3.58	4.30	5.76	8.02	10.72	16.24	26.79
rpma_read() rand	2.86	3.12	3.89	4.64	6.24	8.64	11.58	23.17	35.07
rpma_read() seq	2.96	3.34	3.82	4.64	6.11	8.77	10.94	23.17	35.07

Figure 3. Latency (lat\_pctl\_99.99): ib\_read\_lat vs rpma\_read() from DRAM



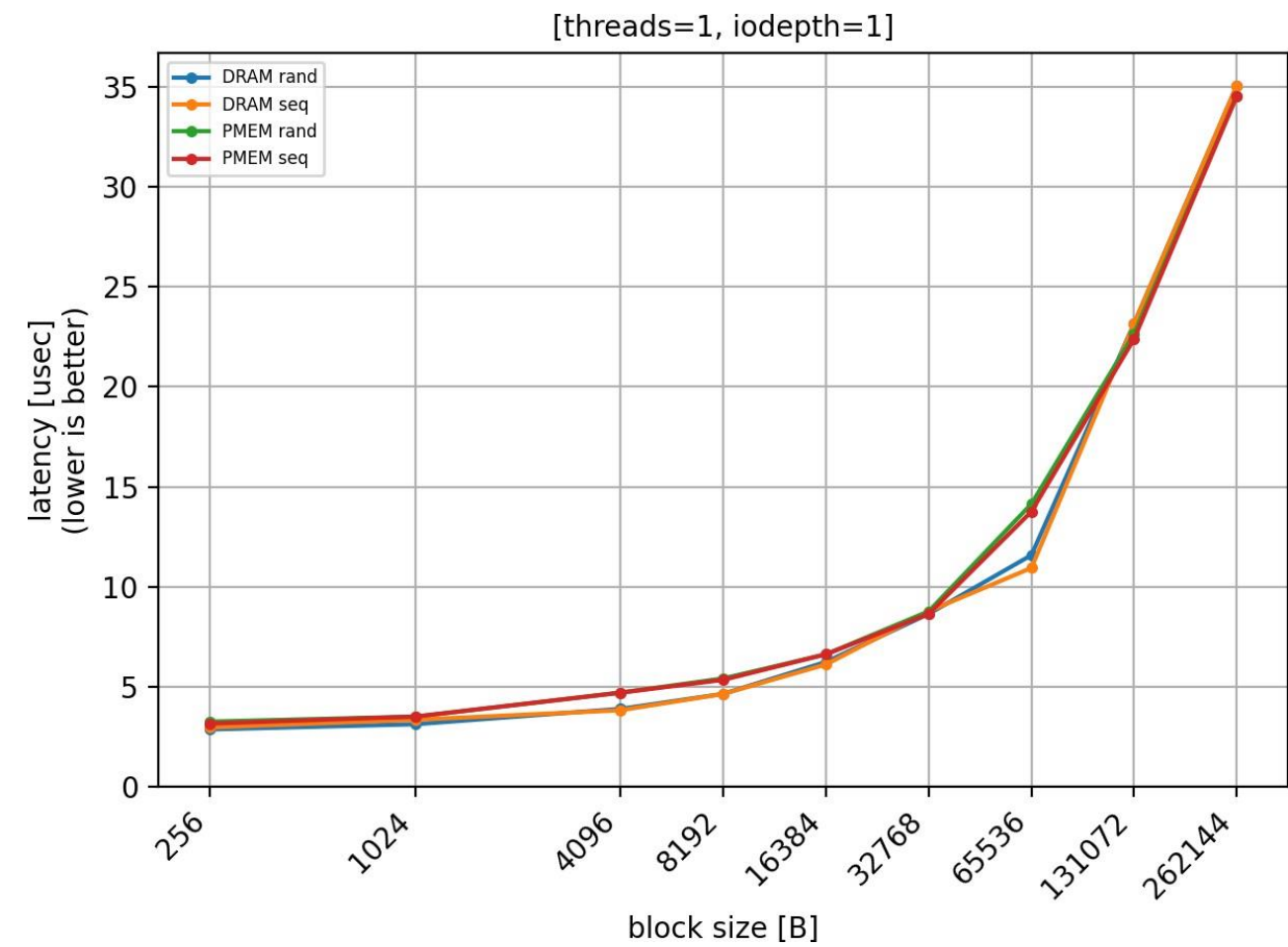
	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read() rand	3.95	4.19	5.54	5.79	7.52	9.92	12.86	24.45	36.10
rpma_read() seq	3.95	4.26	4.96	5.73	7.26	10.18	12.35	24.19	36.10

Figure 4. Latency (lat\_avg): rpma\_read() from DRAM vs from PMEM



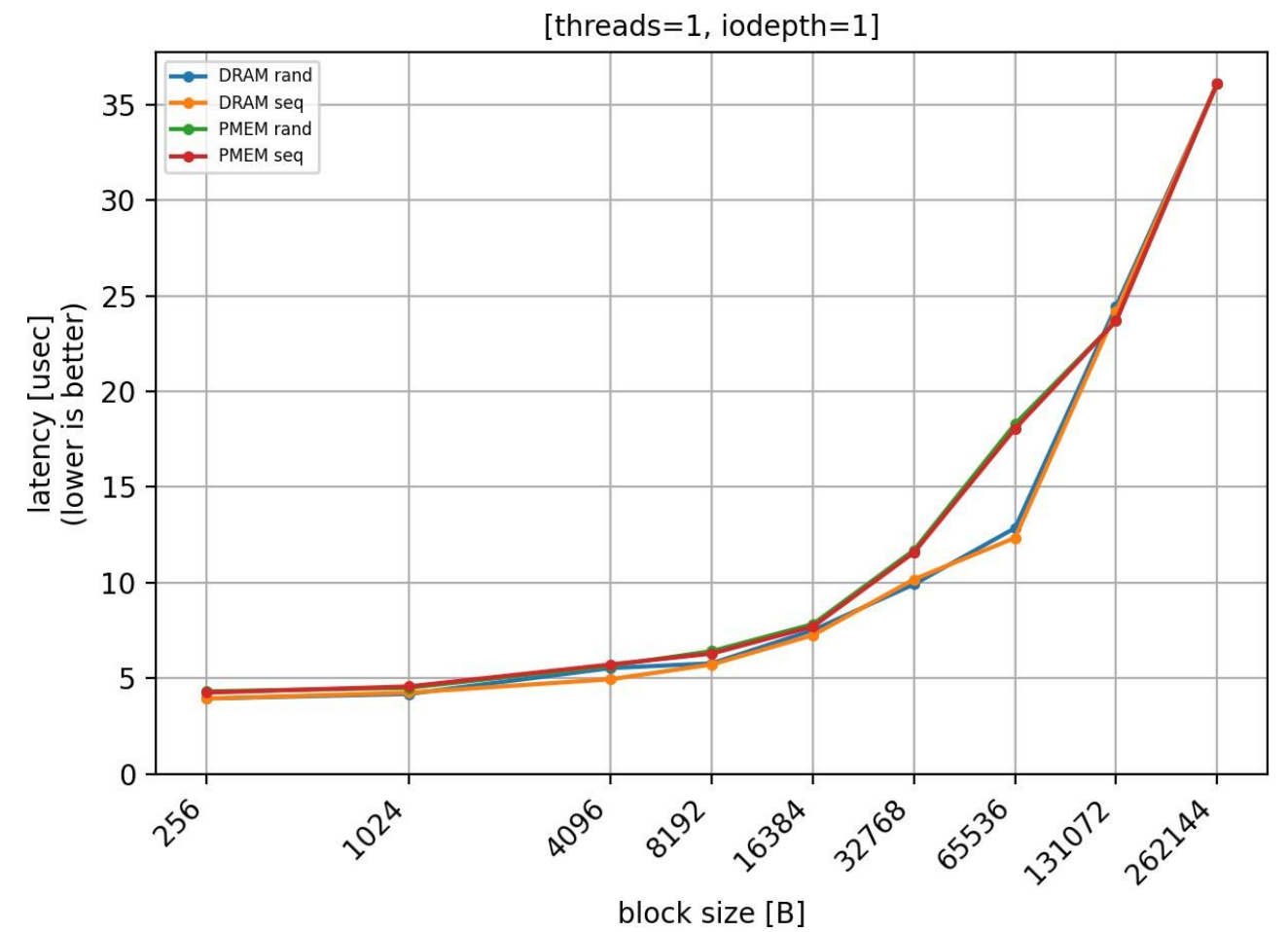
	256	1024	4096	8192	16384	32768	65536	131072	262144
DRAM rand	2.03	2.23	2.97	3.79	5.40	7.82	10.53	15.86	26.75
DRAM seq	1.99	2.17	2.97	3.69	5.29	7.82	10.51	15.85	26.78
PMEM rand	2.29	2.48	3.41	3.99	5.19	7.17	9.83	15.20	26.62
PMEM seq	2.26	2.48	3.40	3.99	5.17	7.08	9.73	15.12	26.54

Figure 5. Latency (lat\_pctl\_99.9): rpma\_read() from DRAM vs from PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
DRAM rand	2.86	3.12	3.89	4.64	6.24	8.64	11.58	23.17	35.07
DRAM seq	2.96	3.34	3.82	4.64	6.11	8.77	10.94	23.17	35.07
PMEM rand	3.25	3.50	4.70	5.41	6.62	8.77	14.14	22.66	34.56
PMEM seq	3.15	3.50	4.70	5.34	6.62	8.64	13.76	22.40	34.56

Figure 6. Latency (lat\_pctl\_99.99): rpma\_read() from DRAM vs from PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
DRAM rand	3.95	4.19	5.54	5.79	7.52	9.92	12.86	24.45	36.10
DRAM seq	3.95	4.26	4.96	5.73	7.26	10.18	12.35	24.19	36.10
PMEM rand	4.32	4.51	5.66	6.43	7.84	11.71	18.30	23.68	36.10
PMEM seq	4.26	4.58	5.73	6.30	7.71	11.58	18.05	23.68	36.10

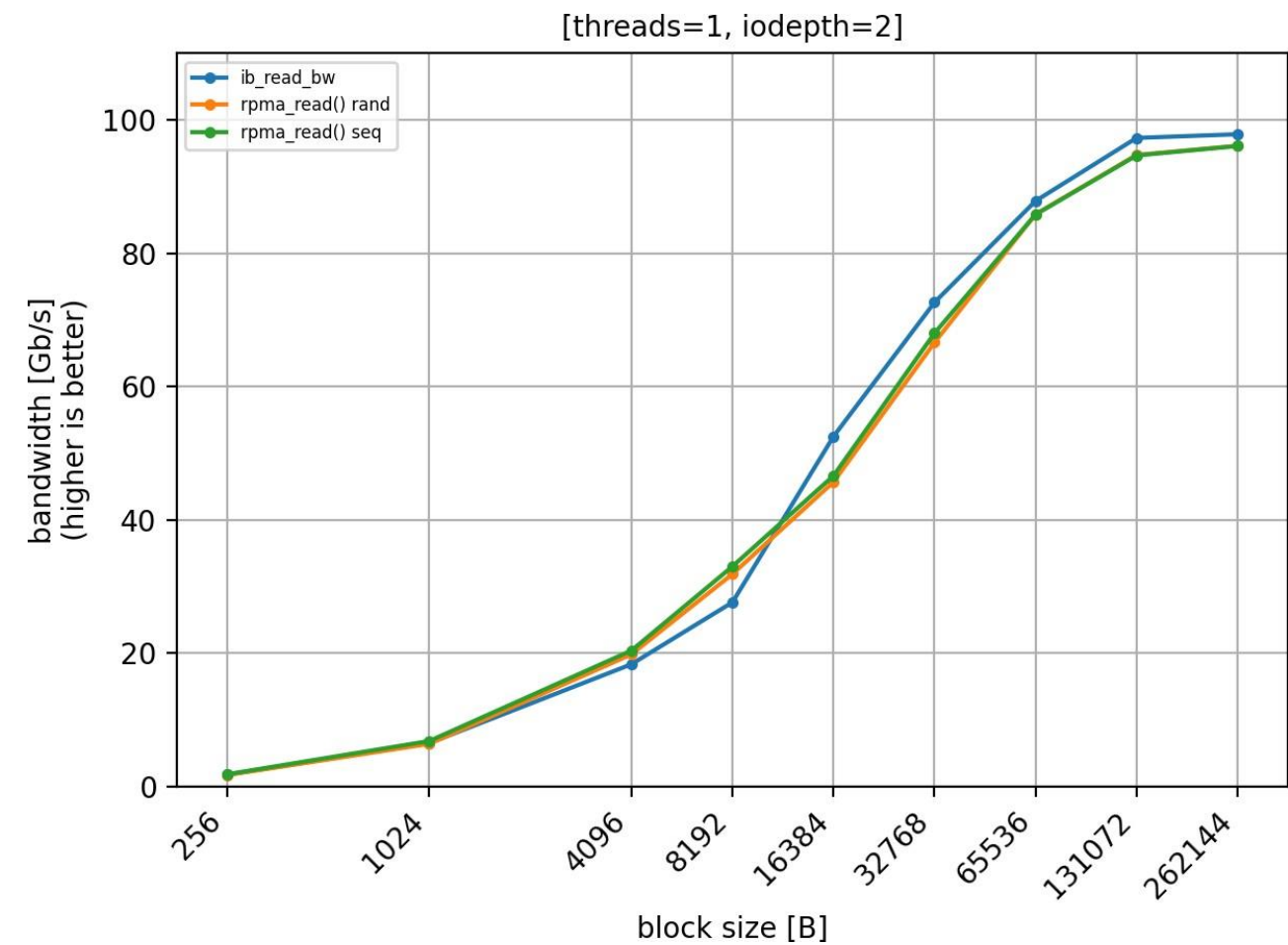
Read from PMem: Bandwidth

Comparing the bandwidth of `rpma_read()` from PMem on **the RPMA Target** to the bandwidth of `rpma_read()` from DRAM on **the RPMA Target** using the `ib_read_bw` as the baseline.

Item	Description
Server - <code>ib_read_bw</code> configuration	<code>--size \$blocksize --qp 1 --tx-depth=2</code> <code>--size 4096 --qp \$threads --tx-depth=2</code>
Client - <code>ib_read_bw</code> configuration	<code>--iters \$iters --size \$blocksize --qp 1 --tx-depth=2 \</code> <code>--report_gbits \$serverip</code> <code>--iters \$iters --size 4096 --qp \$threads --tx-depth=2 \</code> <code>--report_gbits \$serverip</code>
Server - FIO engine configuration	<code>[global]</code> <code>ioengine=librpma_apm_server</code> <code>create_serialize=0</code> <code>kb_base=1000</code> <code>serverip=\$serverip</code> <code>port=7204</code> <code>thread</code>  <code>[server]</code> <code>direct_write_to_pmem={0, 1} # 1 for Device DAX</code> <code>numjobs=1</code> <code>size=100MiB</code> <code>filename={malloc, /dev/dax/path}</code>

Item	Description
Client - FIO engine configuration	<pre>[global] ioengine=librpma_apm_client create_serialize=0 serverip=\$serverip port=7204 thread disable_clat=1 lat_percentiles=1 percentile_list=99.0:99.9:99.99:99.999  [client] numjobs=\$numjobs group_reporting=1 iodepth=2 readwrite={read, randread} blocksize=\$blocksize ramp_time=15 time_based runtime=60</pre>

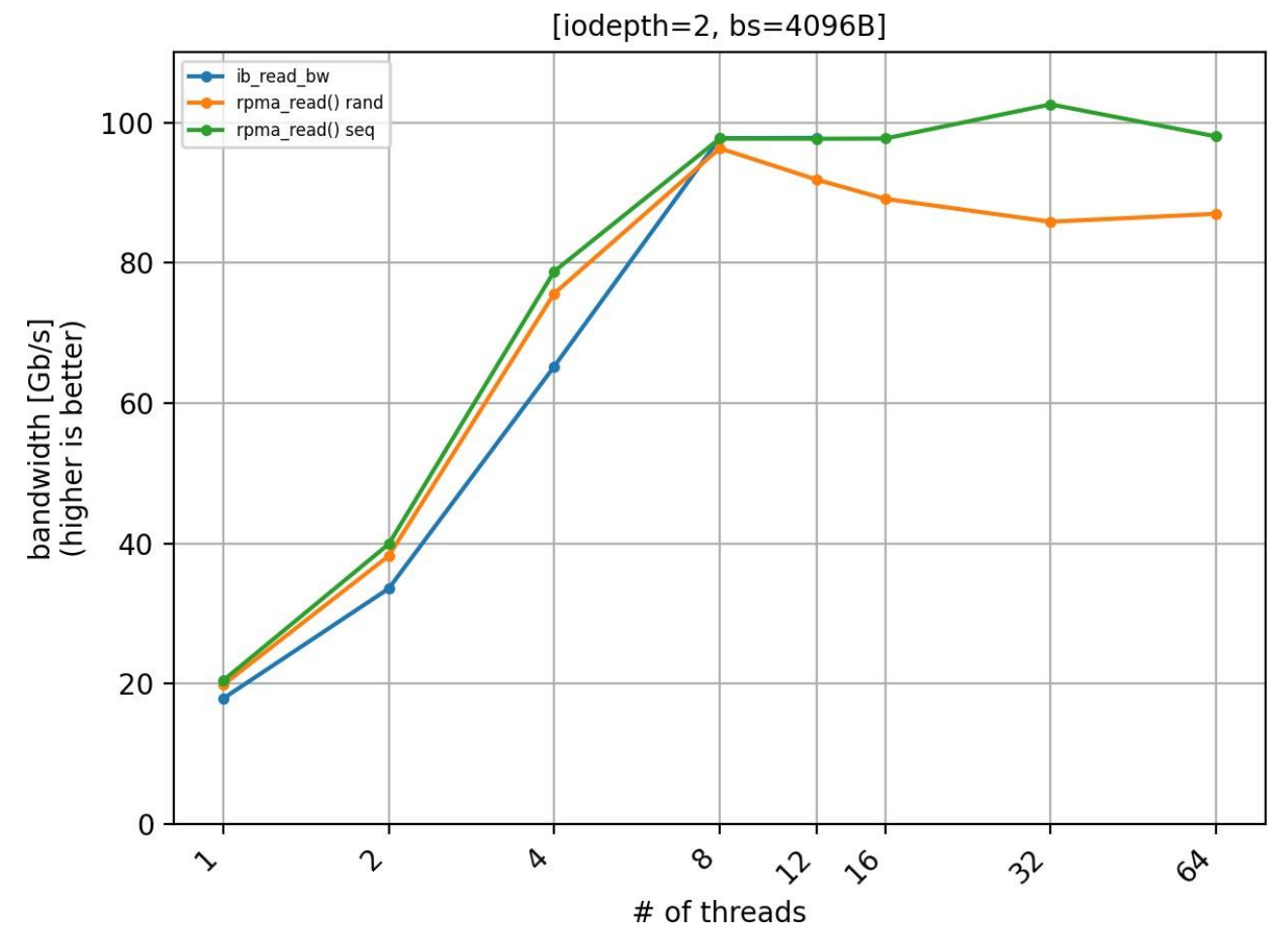
Figure 7. Bandwidth (bs): ib\_read\_bw() vs rpma\_read() from DRAM



	256	1024	4096	8192	16384	32768	65536	131072	262144
ib_read_bw	1.78	6.52	18.35	27.64	52.52	72.64	87.86	97.32	97.87
rpma_read() rand	1.73	6.41	19.83	31.88	45.69	66.62	85.87	94.76	96.13
rpma_read() seq	1.84	6.83	20.40	33.03	46.62	68.02	85.83	94.69	96.11

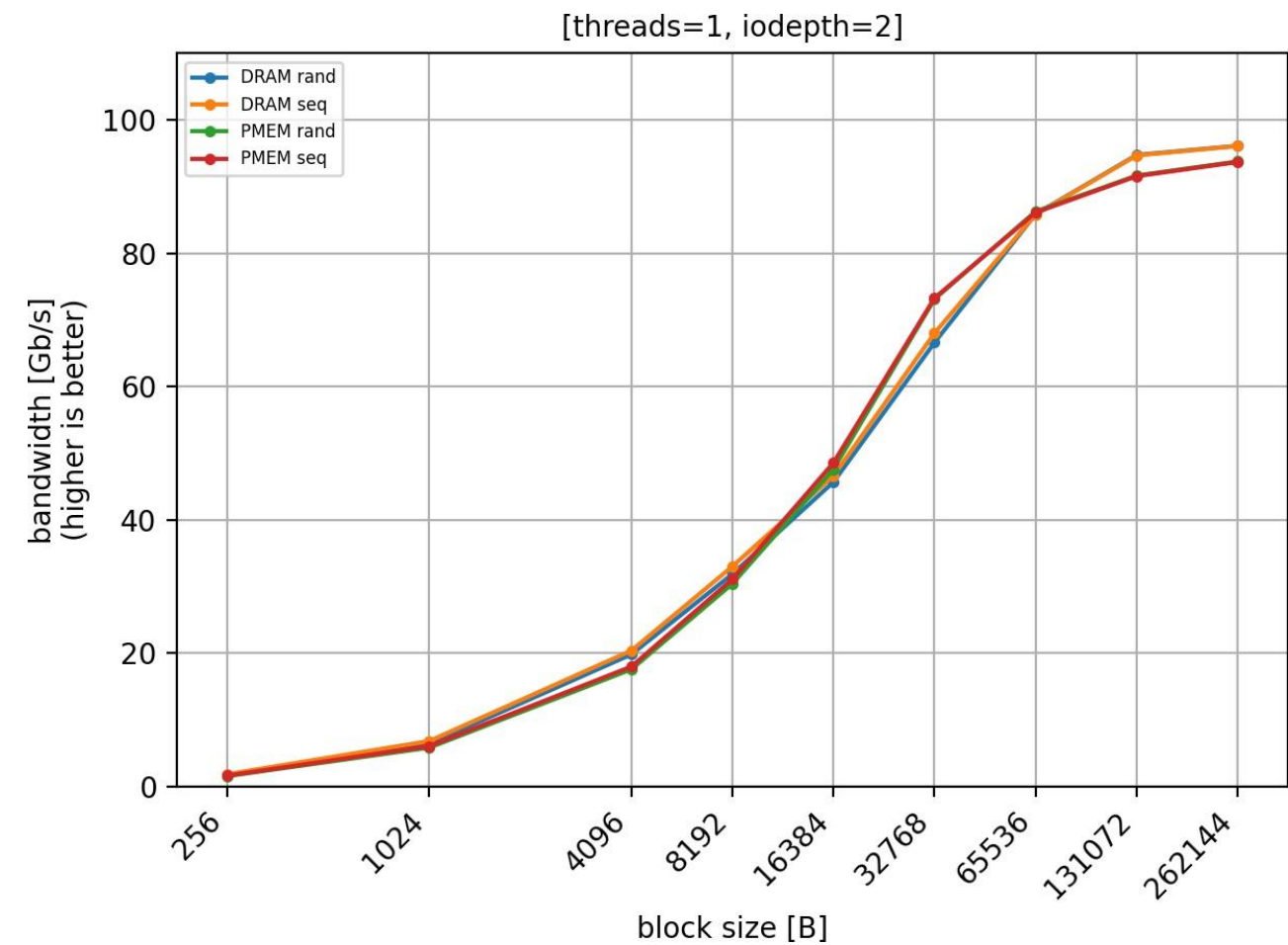


Figure 8. Bandwidth (threads): ib\_read\_bw() vs rpma\_read() from DRAM



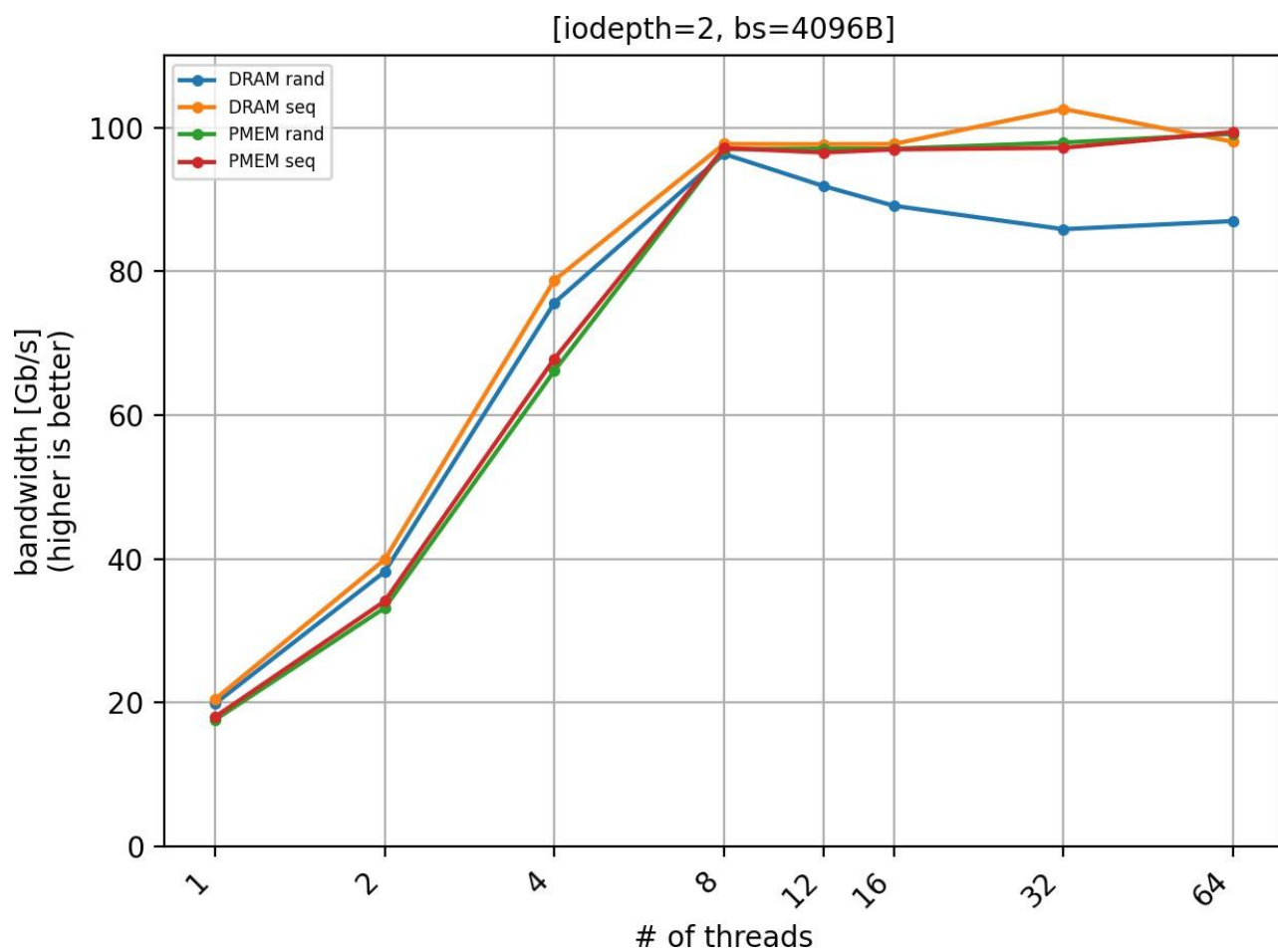
	1	2	4	8	12	16	32	64
ib_read_bw	17.89	33.57	65.24	97.78	97.82	-	-	-
rpma_read() rand	19.85	38.22	75.64	96.34	91.88	89.13	85.87	87.00
rpma_read() seq	20.46	39.91	78.79	97.73	97.70	97.73	102.60	98.03

Figure 9. Bandwidth (bs): rpma\_read() from DRAM vs from PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
DRAM rand	1.73	6.41	19.83	31.88	45.69	66.62	85.87	94.76	96.13
DRAM seq	1.84	6.83	20.40	33.03	46.62	68.02	85.83	94.69	96.11
PMEM rand	1.57	5.84	17.59	30.44	47.61	73.13	86.28	91.67	93.79
PMEM seq	1.63	6.07	18.00	31.10	48.58	73.29	86.13	91.58	93.70

Figure 10. Bandwidth (threads): rpma\_read() from DRAM vs from PMEM



	1	2	4	8	12	16	32	64
DRAM rand	19.85	38.22	75.64	96.34	91.88	89.13	85.87	87.00
DRAM seq	20.46	39.91	78.79	97.73	97.70	97.73	102.60	98.03
PMEM rand	17.60	33.18	66.17	97.00	97.07	97.07	97.93	99.13
PMEM seq	18.00	34.11	67.86	97.16	96.54	96.97	97.19	99.41

## Write to PMem

Benchmarking two ways of writing data persistently to **the RPMA Target**: *Appliance Persistency Method (APM)* and *General Purpose Persistency Method (GPSPM)*. Where:

- ♦ **APM** uses `rpma_flush()` following a sequence of `rpma_write()` operations to provide the remote persistency. This method requires **the RPMA Target** to be capable of *Direct Write to PMem* (for details please see [Direct Write to PMem](#)).
- ♦ **GPSPM** uses `rpma_send()` and `rpma_recv()` operations for sending requests to **the RPMA Target** to assure persistency of the data written using `rpma_write()`. **The RPMA Target** has to provide a thread handling these requests e.g. persisting the data using the `pmem_persist()` operation. When the persistency of the data is assured the response is sent back using also `rpma_send()` and `rpma_recv()`. Depending on how the thread polls for incoming requests you may distinguish two modes:
  - **GPSPM-RT** where the thread polling for incoming requests busy-wait for them (`busy_wait_polling=1`) and
  - **GPSPM** where the thread schedule to be wakened up when a request will appear (`busy_wait_polling=0`). Picking one of these polling modes over another introduces specific challenges and benefits.

For more details on **APM** and **GPSPM** please see: "[Persistent Memory Replication Over Traditional RDMA Part 1: Understanding Remote Persistent Memory](#)" Chapter "Two Remote Replication Methods".

As a baseline is used **APM** to DRAM on **the RPMA Target** with *Direct Write to PMem* capability **disabled**. Such a configuration allows benefiting from available caching on **the RPMA Target**. Note that data, in this case, is not written persistently on **the RPMA Target** but this configuration is used to show the limit of what is possible regarding data transmission in general and how big is potential overhead when transmitting data using **APM** or **GPSPM** while both of these provide the remote persistency additionally.

### Write to PMem: Latency

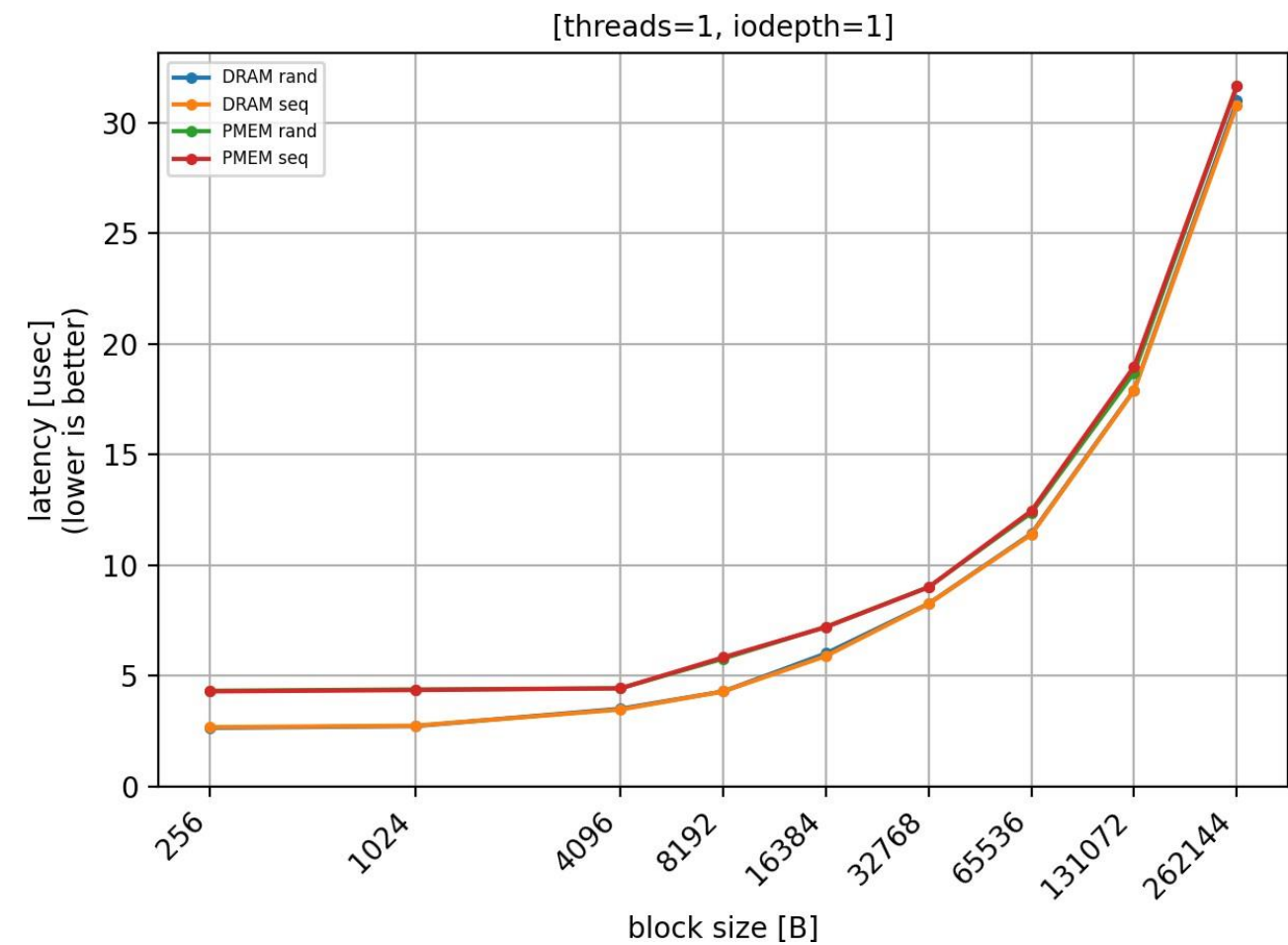
Comparing the latency of **APM** to PMem on **the RPMA Target** (with *Direct Write to PMem*) vs the latency of **APM** to DRAM on **the RPMA Target** (with *Direct Write to PMem* disabled) (as a baseline) vs the latency of **GPSPM(-RT)** to PMem on **the RPMA Target** (with *Direct Write to PMem* disabled).

Item	Description
APM Server - FIO engine configuration	<pre>[global] ioengine=librpma_apm_server create_serialize=0 kb_base=1000 serverip=\$serverip port=7204 thread  [server] direct_write_to_pmem={0, 1} # 1 for Device DAX numjobs=1 size=100MiB filename={malloc, /dev/dax/path}</pre>

Item	Description
APM Client - FIO engine configuration	<pre> [global] ioengine=librpma_apm_client create_serialize=0 serverip=\$serverip port=7204 thread disable_clat=1 lat_percentiles=1 percentile_list=99.0:99.9:99.99:99.999  [client] sync=1 readwrite={write, randwrite} blocksize=\$blocksize ramp_time=15 time_based runtime=60 </pre>
GPSPM(-RT) Server - FIO engine configuration	<pre> [global] ioengine=librpma_gpssp_server create_serialize=0 kb_base=1000 serverip=\$serverip port=7204 thread  [server] direct_write_to_pmem=0 numjobs=1 iodepth=1 size=100MiB filename=/dev/dax/path busy_wait_polling={0, 1} # 1 for GPSPM-RT  time_based runtime=365d </pre>

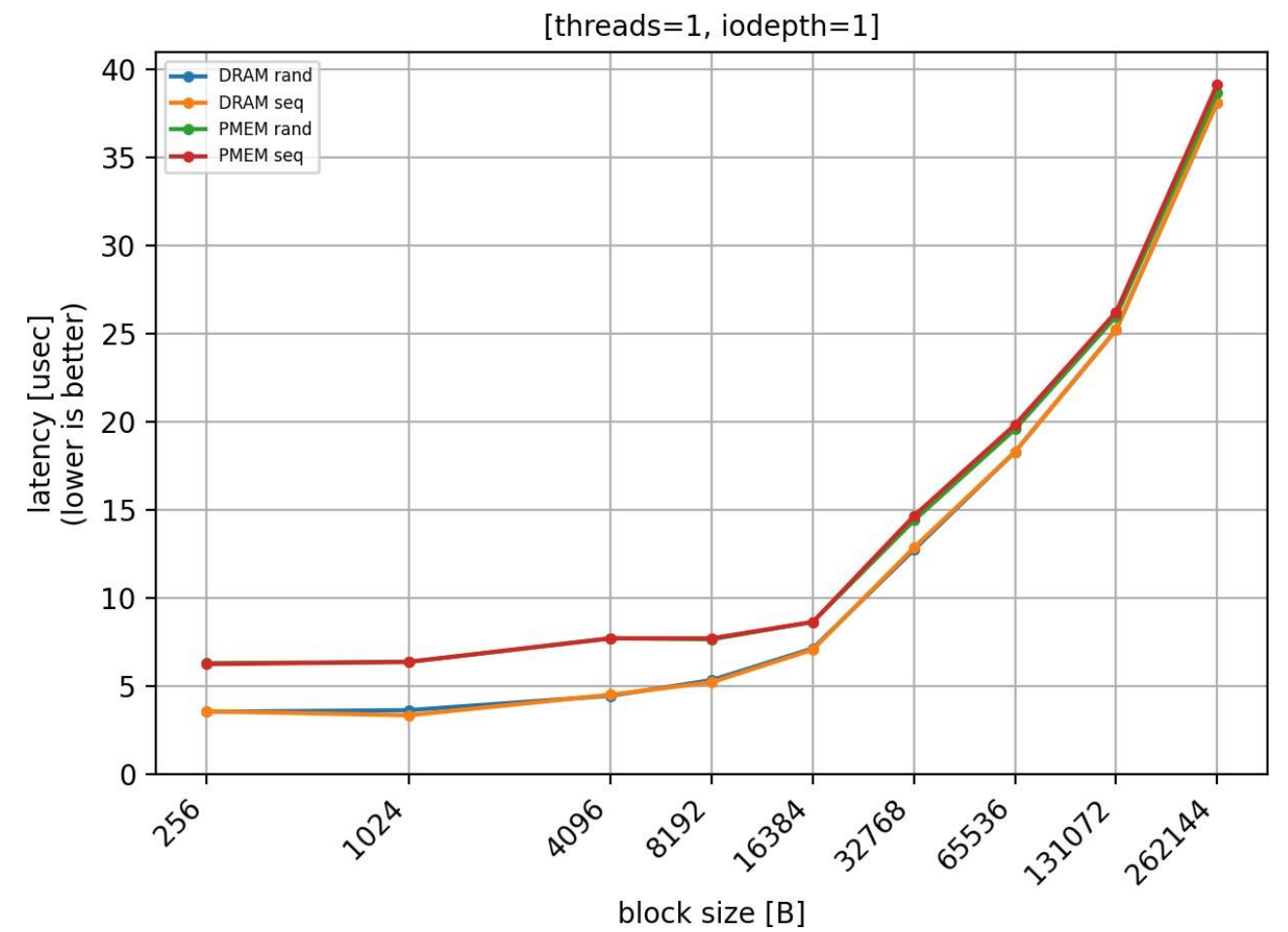
Item	Description
GPSPM(-RT) Client - FIO engine configuration	<pre>[global] ioengine=librpma_gpspm_client create_serialize=0 serverip=\$serverip port=7204 thread disable_clat=1 lat_percentiles=1 percentile_list=99.0:99.9:99.99:99.999  [client] sync=1 readwrite={write, randwrite} blocksize=\$blocksize ramp_time=15 time_based runtime=60</pre>

Figure 11. Latency (lat\_avg): APM to DRAM (DDIO=ON) vs PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
DRAM rand	2.65	2.73	3.52	4.30	6.02	8.27	11.43	17.92	30.99
DRAM seq	2.68	2.75	3.48	4.30	5.90	8.26	11.40	17.90	30.80
PMEM rand	4.33	4.39	4.44	5.78	7.22	9.02	12.37	18.70	31.68
PMEM seq	4.30	4.36	4.44	5.84	7.21	9.01	12.46	18.98	31.65

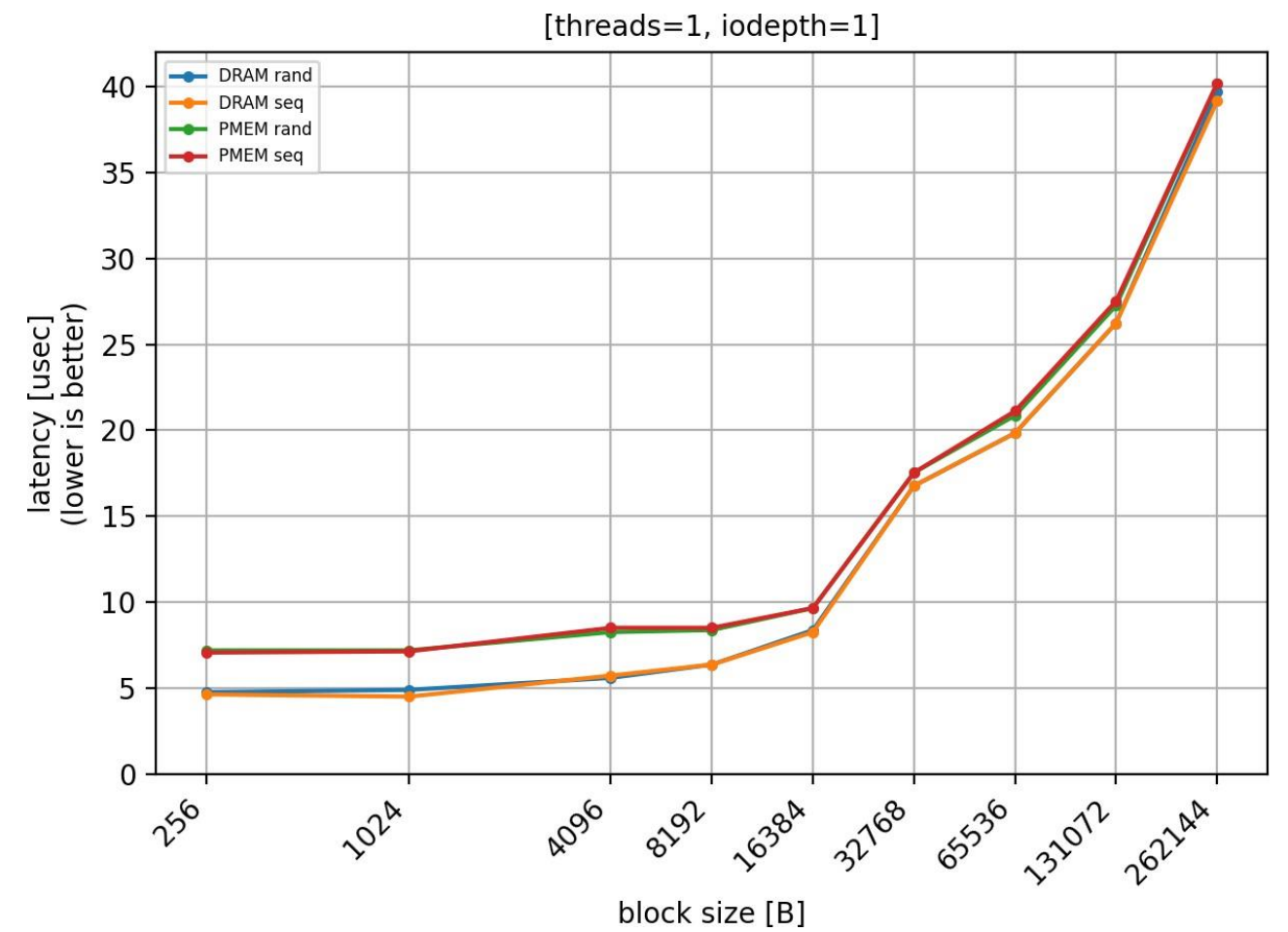
Figure 12. Latency (lat\_pctl\_99.9): APM to DRAM (DDIO=ON) vs PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
DRAM rand	3.54	3.63	4.45	5.34	7.14	12.74	18.30	25.22	38.14
DRAM seq	3.57	3.34	4.51	5.22	7.07	12.86	18.30	25.22	38.14
PMEM rand	6.30	6.37	7.71	7.65	8.64	14.40	19.58	25.98	38.66
PMEM seq	6.24	6.37	7.71	7.71	8.64	14.66	19.84	26.24	39.17

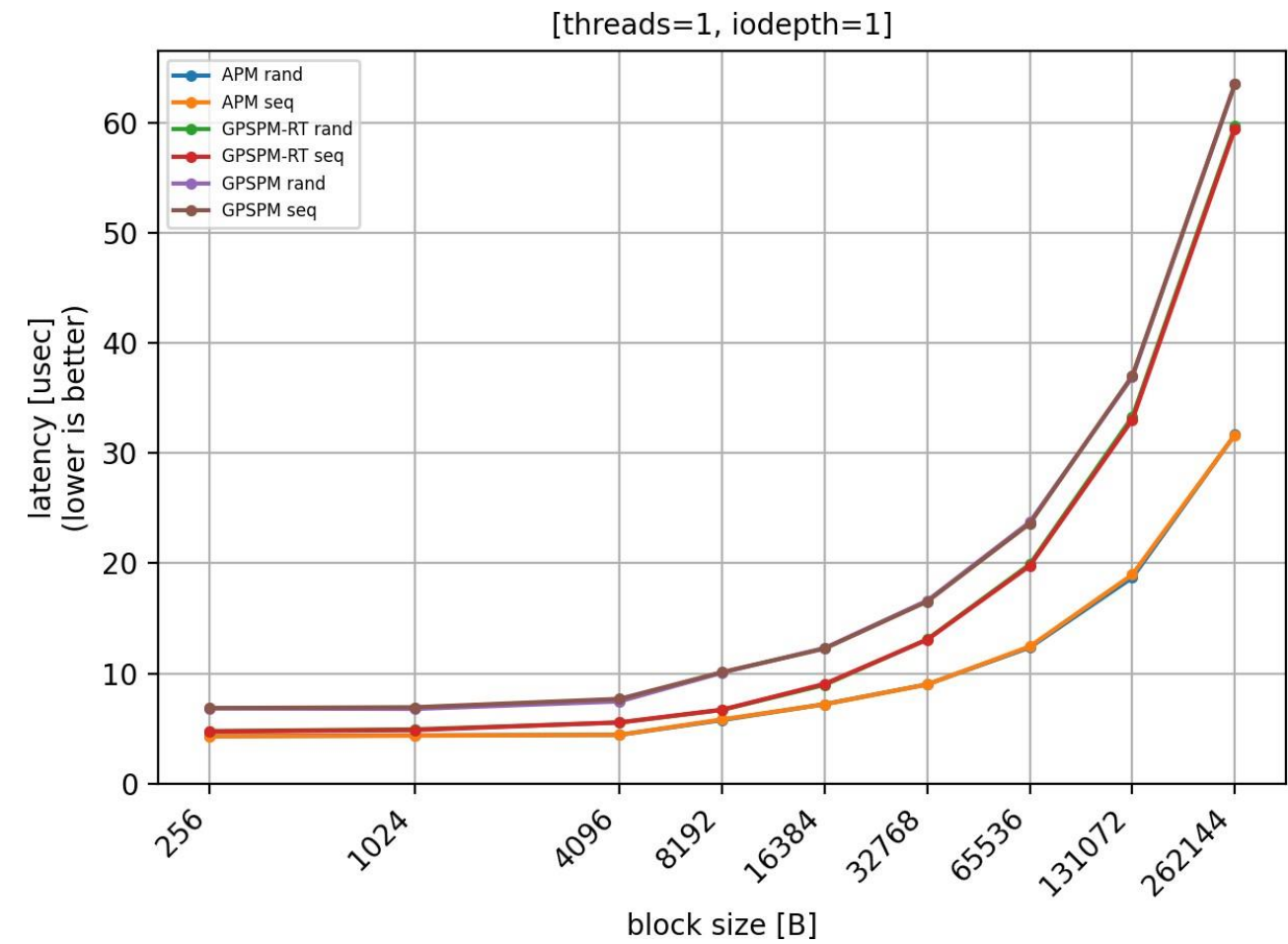


Figure 13. Latency (lat\_pctl\_99.99): APM to DRAM (DDIO=ON) vs PMEM



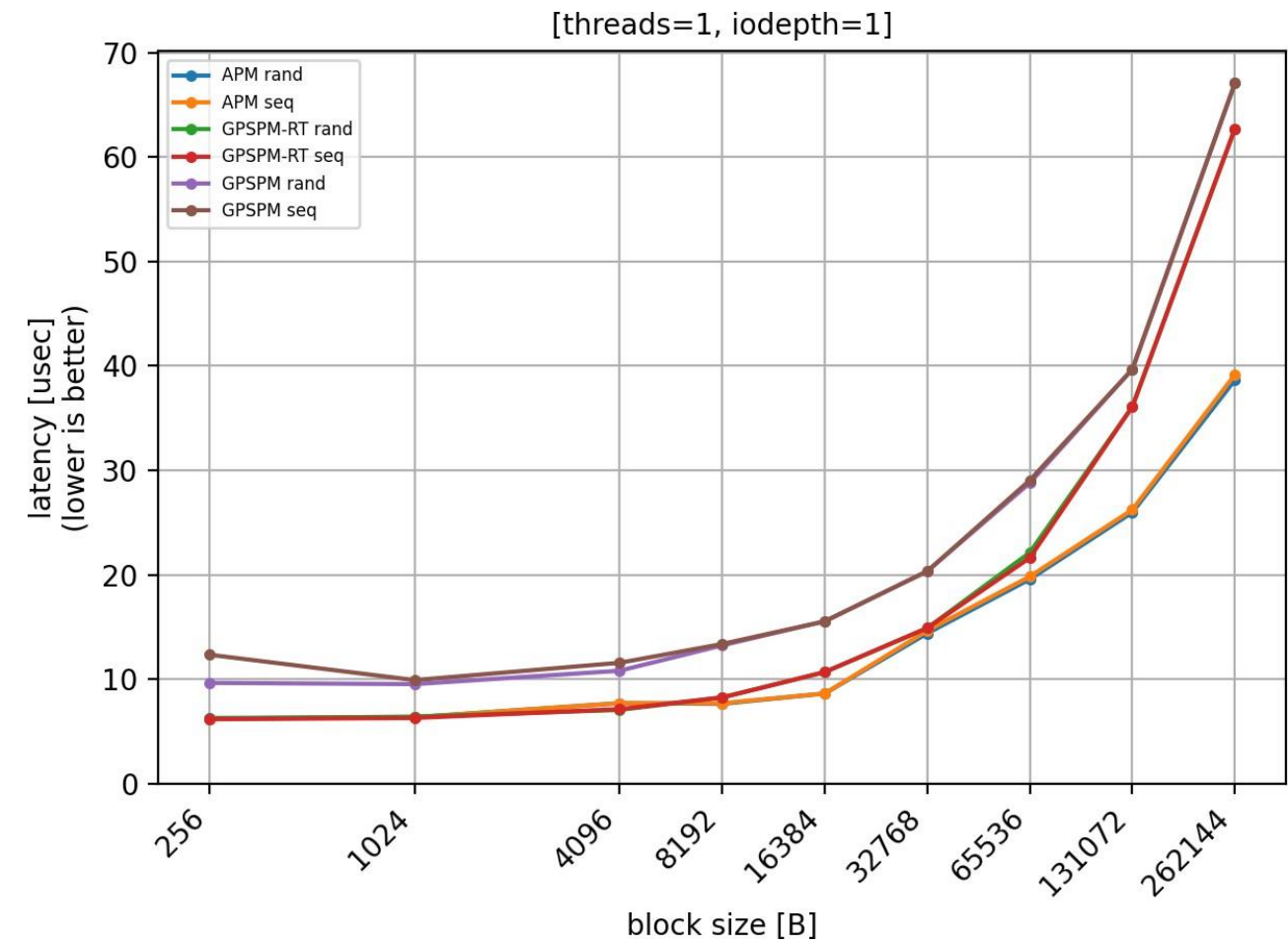
	256	1024	4096	8192	16384	32768	65536	131072	262144
DRAM rand	4.77	4.90	5.60	6.37	8.38	16.77	19.84	26.24	39.68
DRAM seq	4.64	4.51	5.73	6.37	8.26	16.77	19.84	26.24	39.17
PMEM rand	7.20	7.20	8.26	8.38	9.66	17.54	20.86	27.26	40.19
PMEM seq	7.07	7.14	8.51	8.51	9.66	17.54	21.12	27.52	40.19

Figure 14. Latency (lat\_avg): APM to PMEM vs GPSPM(-RT) to PMEM



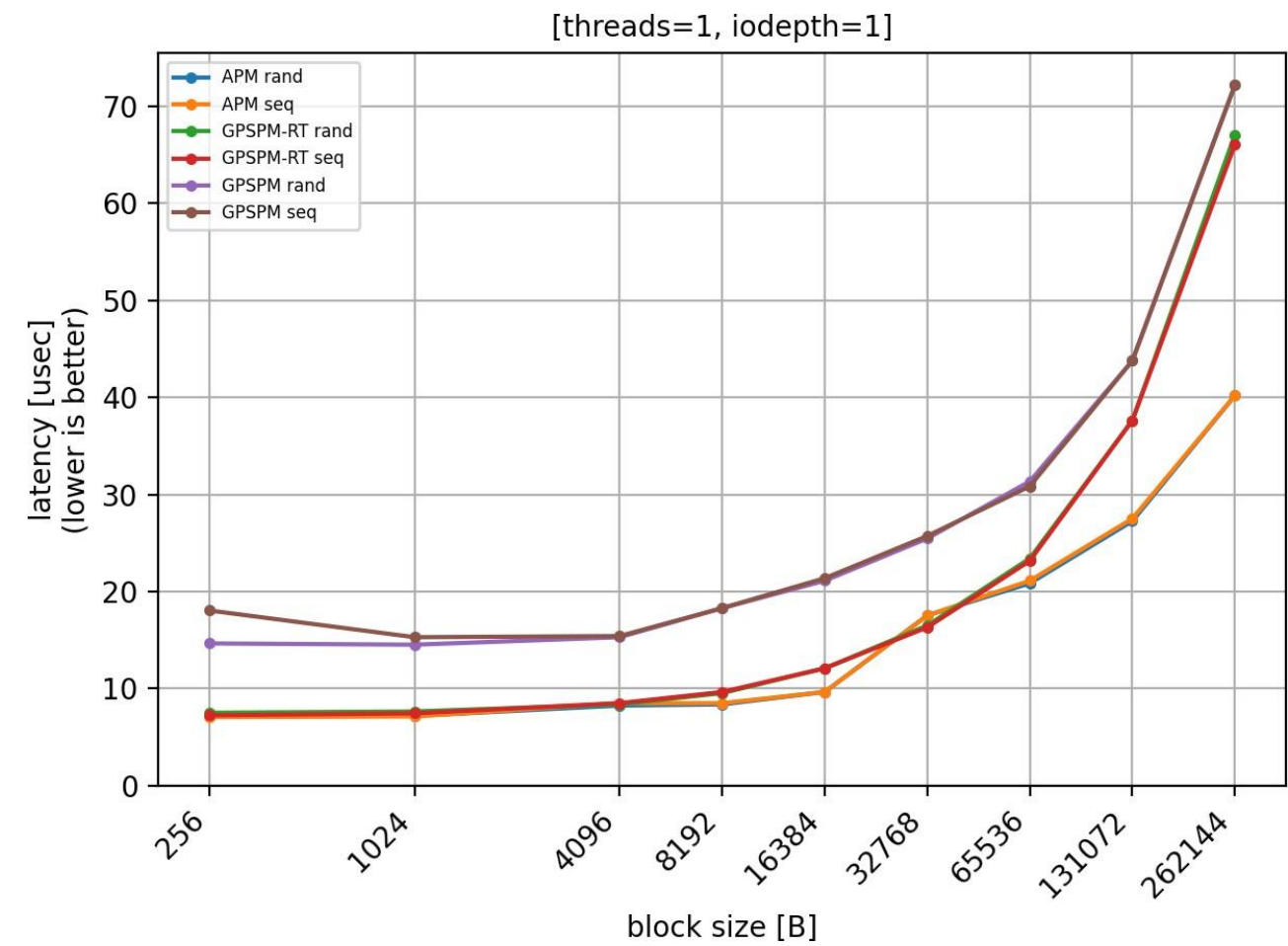
	256	1024	4096	8192	16384	32768	65536	131072	262144
APM rand	4.33	4.39	4.44	5.78	7.22	9.02	12.37	18.70	31.68
APM seq	4.30	4.36	4.44	5.84	7.21	9.01	12.46	18.98	31.65
GPSPM-RT rand	4.78	4.94	5.55	6.70	8.97	13.11	19.92	33.30	59.67
GPSPM-RT seq	4.74	4.87	5.57	6.72	9.07	13.08	19.77	33.02	59.42
GPSPM rand	6.84	6.79	7.47	10.08	12.32	16.62	23.75	36.90	63.49
GPSPM seq	6.87	6.94	7.70	10.13	12.27	16.54	23.61	37.02	63.49

Figure 15. Latency (lat\_pctl\_99.9): APM to PMEM vs GPSPM(-RT) to PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
APM rand	6.30	6.37	7.71	7.65	8.64	14.40	19.58	25.98	38.66
APM seq	6.24	6.37	7.71	7.71	8.64	14.66	19.84	26.24	39.17
GPSPM-RT rand	6.24	6.43	7.07	8.26	10.69	14.91	22.14	36.10	62.72
GPSPM-RT seq	6.18	6.30	7.14	8.26	10.69	14.91	21.63	36.10	62.72
GPSPM rand	9.66	9.54	10.82	13.25	15.55	20.35	28.80	39.68	67.07
GPSPM seq	12.35	9.92	11.58	13.38	15.55	20.35	29.06	39.68	67.07

Figure 16. Latency (lat\_pctl\_99.99): APM to PMEM vs GPSPM(-RT) to PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
APM rand	7.20	7.20	8.26	8.38	9.66	17.54	20.86	27.26	40.19
APM seq	7.07	7.14	8.51	8.51	9.66	17.54	21.12	27.52	40.19
GPSPM-RT rand	7.52	7.65	8.38	9.54	12.10	16.51	23.42	37.63	67.07
GPSPM-RT seq	7.26	7.46	8.51	9.66	12.10	16.32	23.17	37.63	66.05
GPSPM rand	14.66	14.53	15.30	18.30	21.12	25.47	31.36	43.78	72.19
GPSPM seq	18.05	15.30	15.42	18.30	21.38	25.73	30.85	43.78	72.19

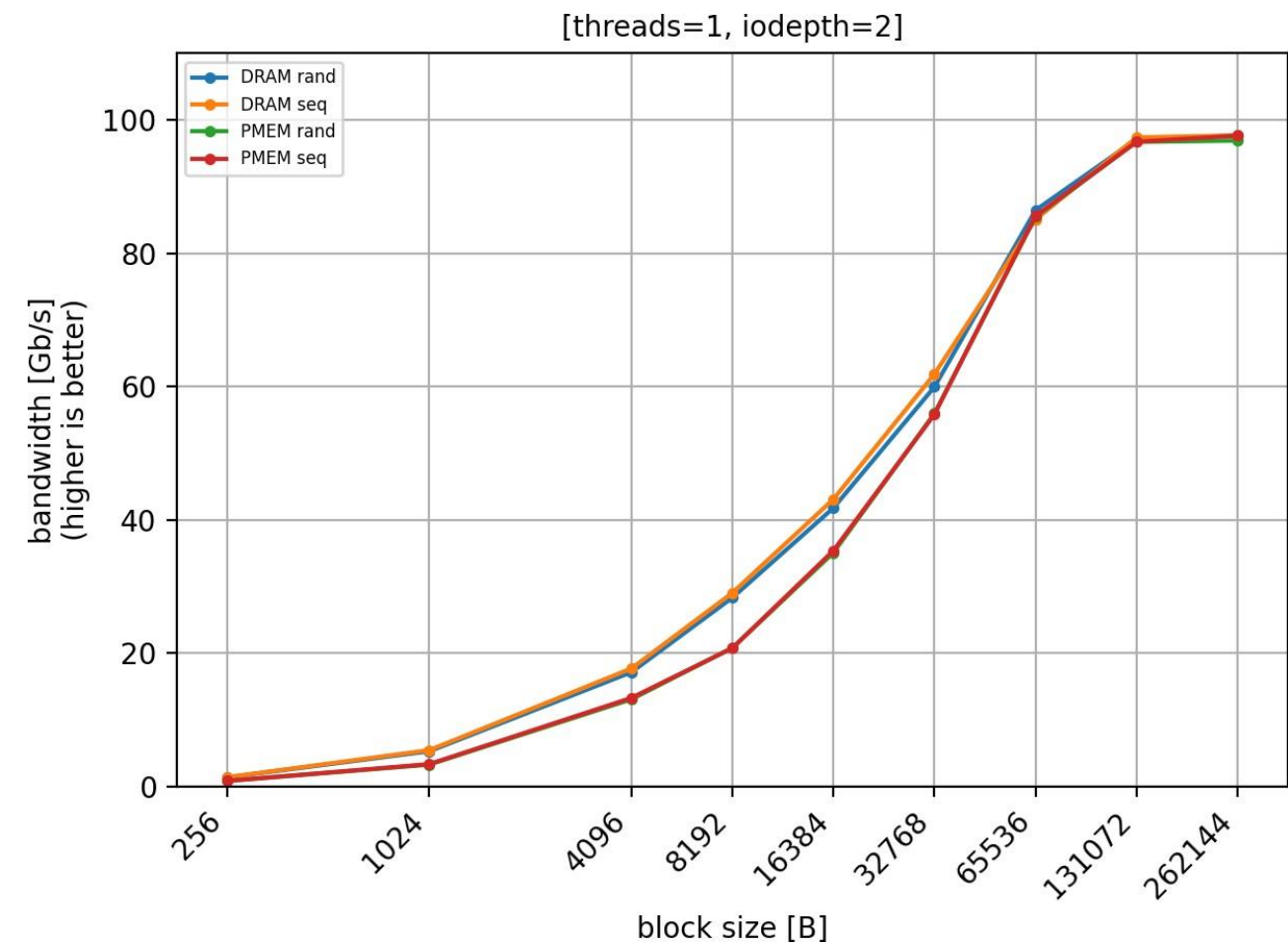
### Write to PMem: Bandwidth

Comparing the bandwidth of **APM** to PMem on **the RPMA Target** (with *Direct Write to PMem*) vs the bandwidth of **APM** to DRAM on **the RPMA Target** (with *Direct Write to PMem* disabled) (as a baseline) vs the bandwidth of **GPSPM(-RT)** to PMem on **the RPMA Target** (with *Direct Write to PMem* disabled).

Item	Description
APM Server - FIO engine configuration	<pre>[global] ioengine=librpma_apm_server create_serialize=0 kb_base=1000 serverip=\$serverip port=7204 thread  [server] direct_write_to_pmem=1 numjobs=\$numjobs size=100MiB filename={malloc, /dev/dax/path}</pre>
APM Client - FIO engine configuration	<pre>[global] ioengine=librpma_apm_client create_serialize=0 serverip=\$serverip port=7204 thread disable_clat=1 lat_percentiles=1 percentile_list=99.0:99.9:99.99:99.999  [client] numjobs=\$numjobs group_reporting=1 iodepth=2 readwrite={write, randwrite} blocksize=\$blocksize ramp_time=15 time_based runtime=60</pre>

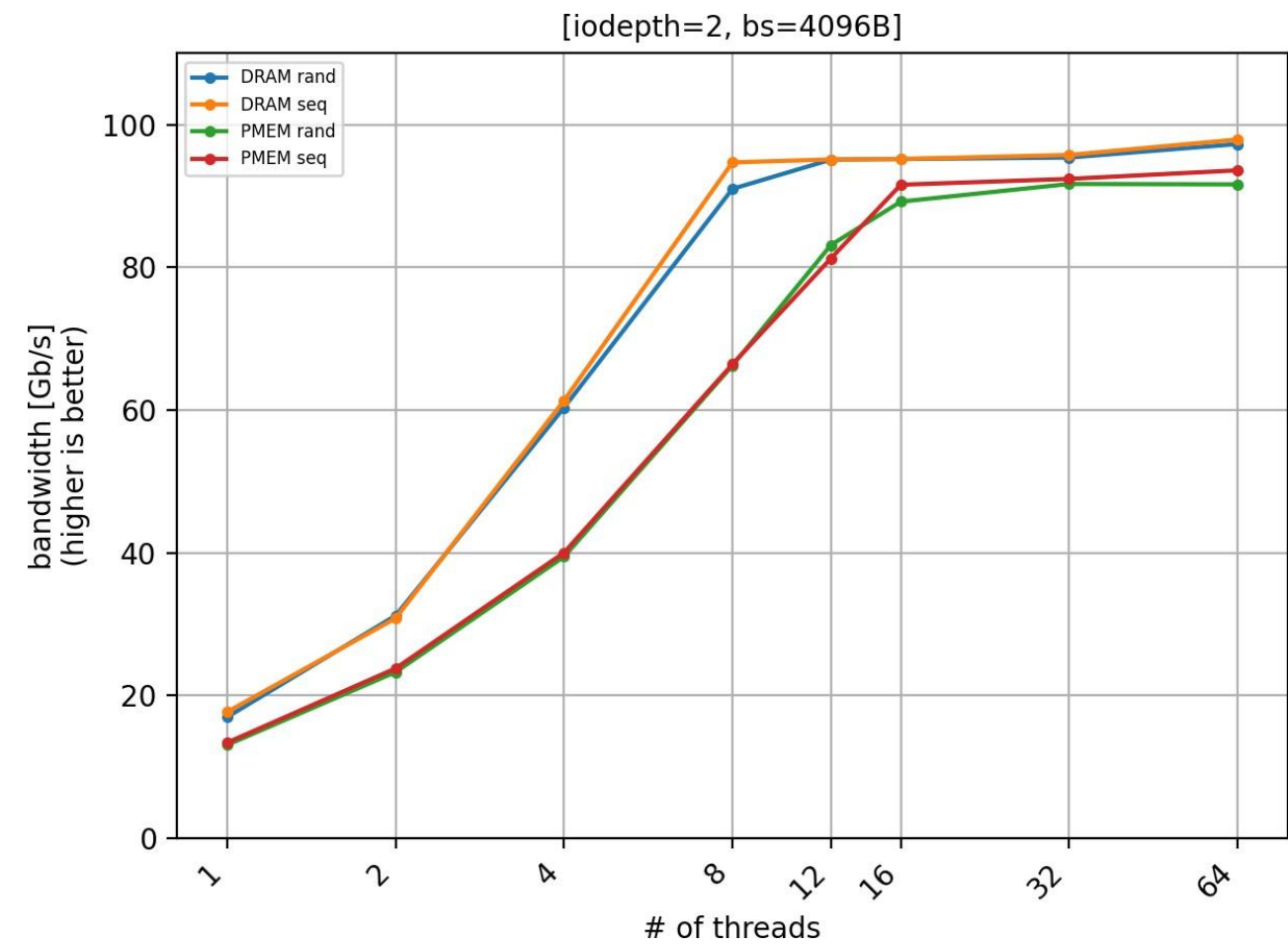
Item	Description
GPSPM(-RT) Server - FIO engine configuration	<pre> [global] ioengine=librpma_gpssp_server create_serialize=0 kb_base=1000 serverip=\$serverip port=7204 thread  [server] direct_write_to_pmem=0 numjobs=\$numjobs iodepth=2 size=100MiB filename=/dev/dax/path busy_wait_polling={0, 1} # 1 for GPSPM-RT  time_based runtime=365d </pre>
GPSPM(-RT) Client - FIO engine configuration	<pre> [global] ioengine=librpma_gpssp_client create_serialize=0 serverip=\$serverip port=7204 thread disable_clat=1 lat_percentiles=1 percentile_list=99.0:99.9:99.99:99.999  [client] numjobs=\$numjobs group_reporting=1 iodepth=2 readwrite={write, randwrite} blocksize=\$blocksize ramp_time=15 time_based runtime=60 </pre>

Figure 17. Bandwidth (bs): APM to DRAM (DDIO=ON) vs to PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
DRAM rand	1.38	5.29	17.15	28.35	41.81	59.97	86.48	97.12	97.69
DRAM seq	1.40	5.46	17.71	29.07	43.12	61.90	85.08	97.40	97.69
PMEM rand	0.84	3.29	13.09	20.81	35.09	55.98	85.50	96.72	96.92
PMEM seq	0.85	3.36	13.30	20.78	35.41	55.86	85.62	96.75	97.65

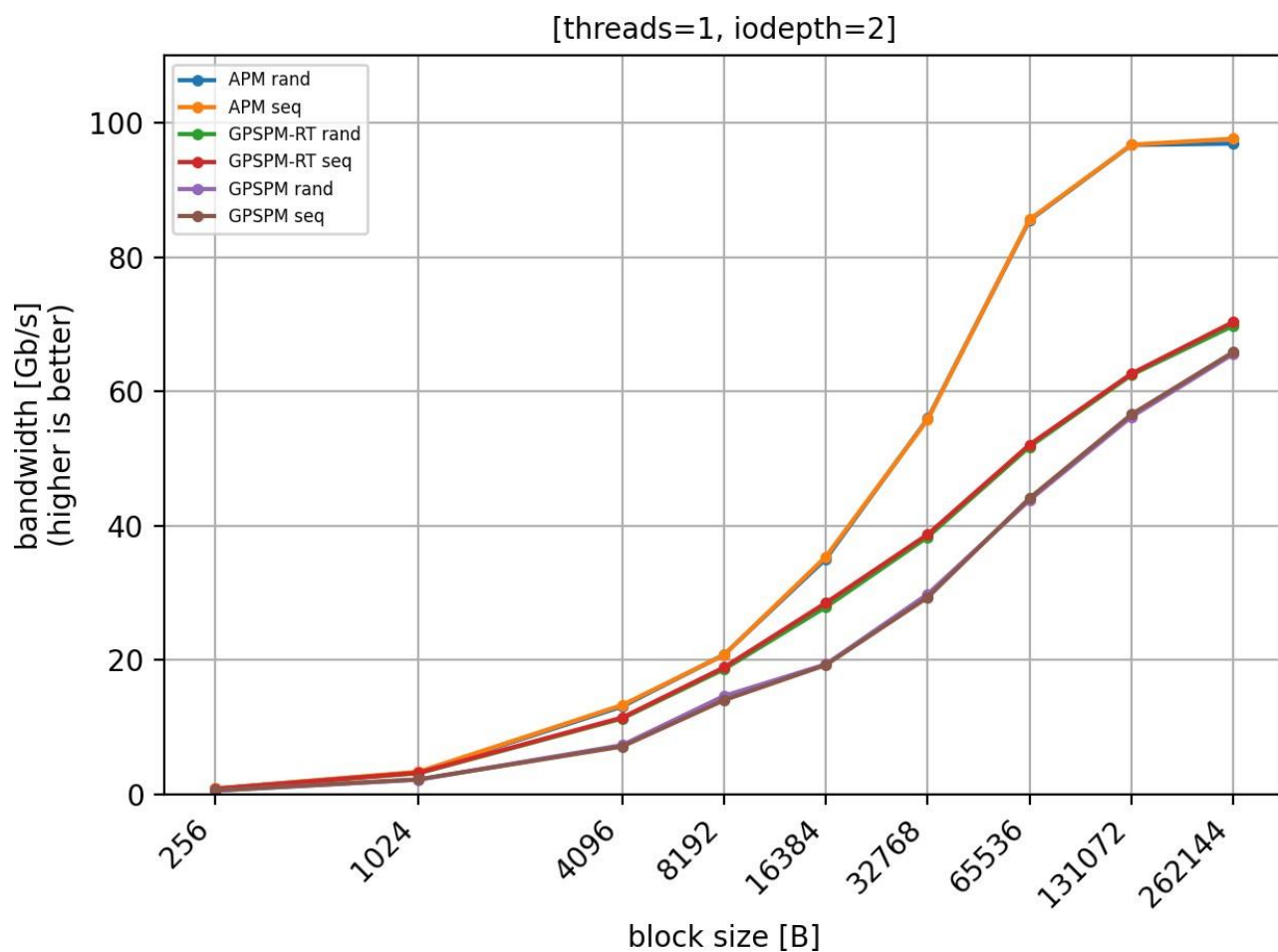
Figure 18. Bandwidth (threads): APM to DRAM (DDIO=ON) vs to PMEM



	1	2	4	8	12	16	32	64
DRAM rand	16.97	31.20	60.34	91.01	95.11	95.19	95.40	97.29
DRAM seq	17.69	30.82	61.34	94.72	95.13	95.19	95.76	97.93
PMEM rand	13.06	23.28	39.45	66.24	83.14	89.22	91.70	91.64
PMEM seq	13.38	23.79	40.01	66.43	81.26	91.59	92.41	93.61

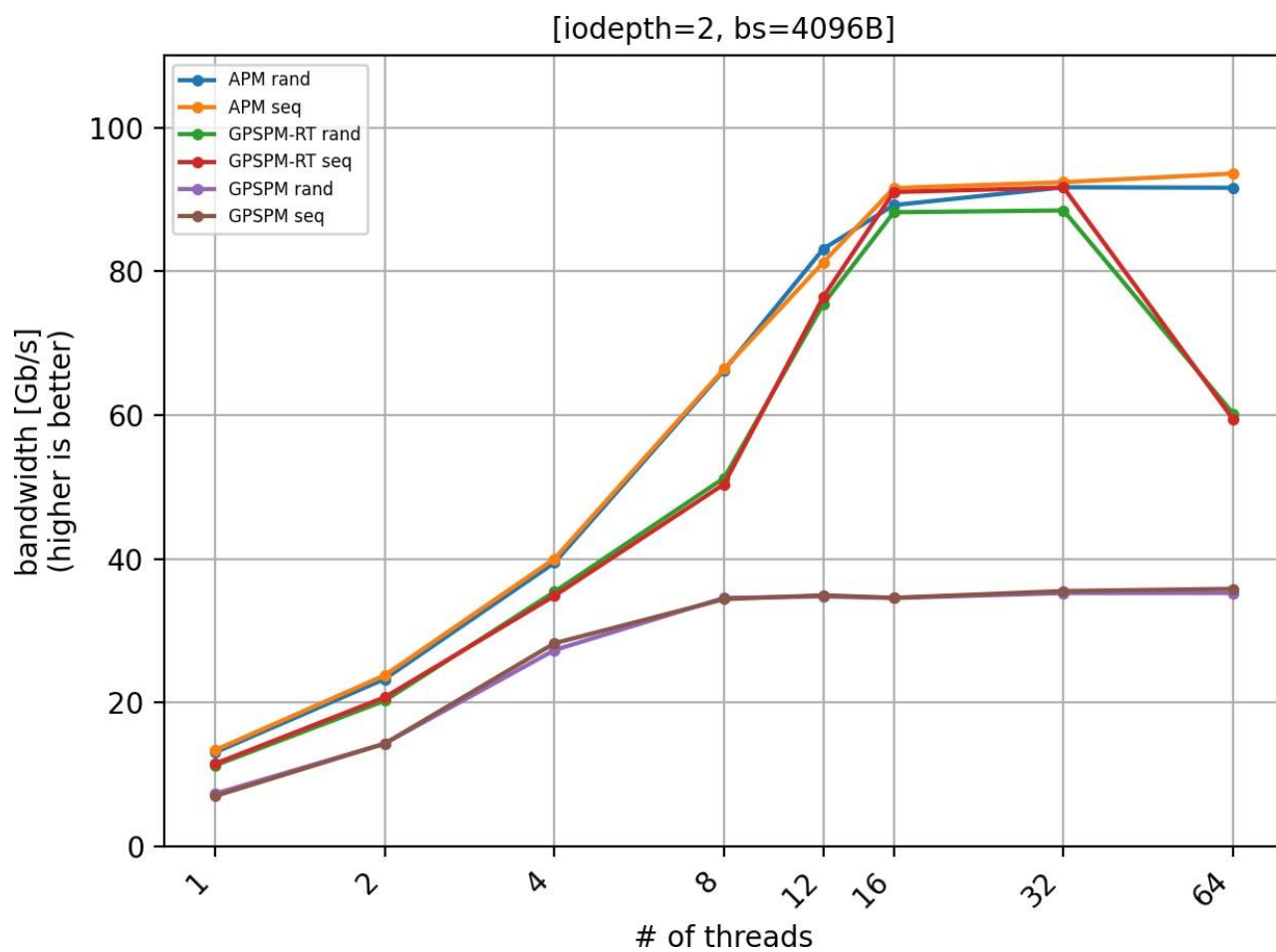


Figure 19. Bandwidth (bs): APM to PMEM vs GPSPM(-RT) to PMEM



	256	1024	4096	8192	16384	32768	65536	131072	262144
APM rand	0.84	3.29	13.09	20.81	35.09	55.98	85.50	96.72	96.92
APM seq	0.85	3.36	13.30	20.78	35.41	55.86	85.62	96.75	97.65
GPSPM-RT rand	0.79	3.15	11.30	18.64	27.91	38.34	51.78	62.46	69.80
GPSPM-RT seq	0.82	3.22	11.44	18.92	28.53	38.73	52.09	62.66	70.35
GPSPM rand	0.56	2.20	7.34	14.66	19.37	29.83	43.83	56.23	65.59
GPSPM seq	0.57	2.23	7.14	14.04	19.29	29.36	44.12	56.58	65.87

Figure 20. Bandwidth (threads): APM to PMEM vs GPSPM(-RT) to PMEM



	1	2	4	8	12	16	32	64
APM rand	13.06	23.28	39.45	66.24	83.14	89.22	91.70	91.64
APM seq	13.38	23.79	40.01	66.43	81.26	91.59	92.41	93.61
GPSPM-RT rand	11.25	20.26	35.41	51.25	75.48	88.23	88.48	60.20
GPSPM-RT seq	11.49	20.73	34.85	50.40	76.46	91.04	91.67	59.40
GPSPM rand	7.32	14.29	27.30	34.57	34.76	34.54	35.24	35.25
GPSPM seq	6.99	14.29	28.25	34.39	34.93	34.58	35.52	35.83

# Mix against PMem

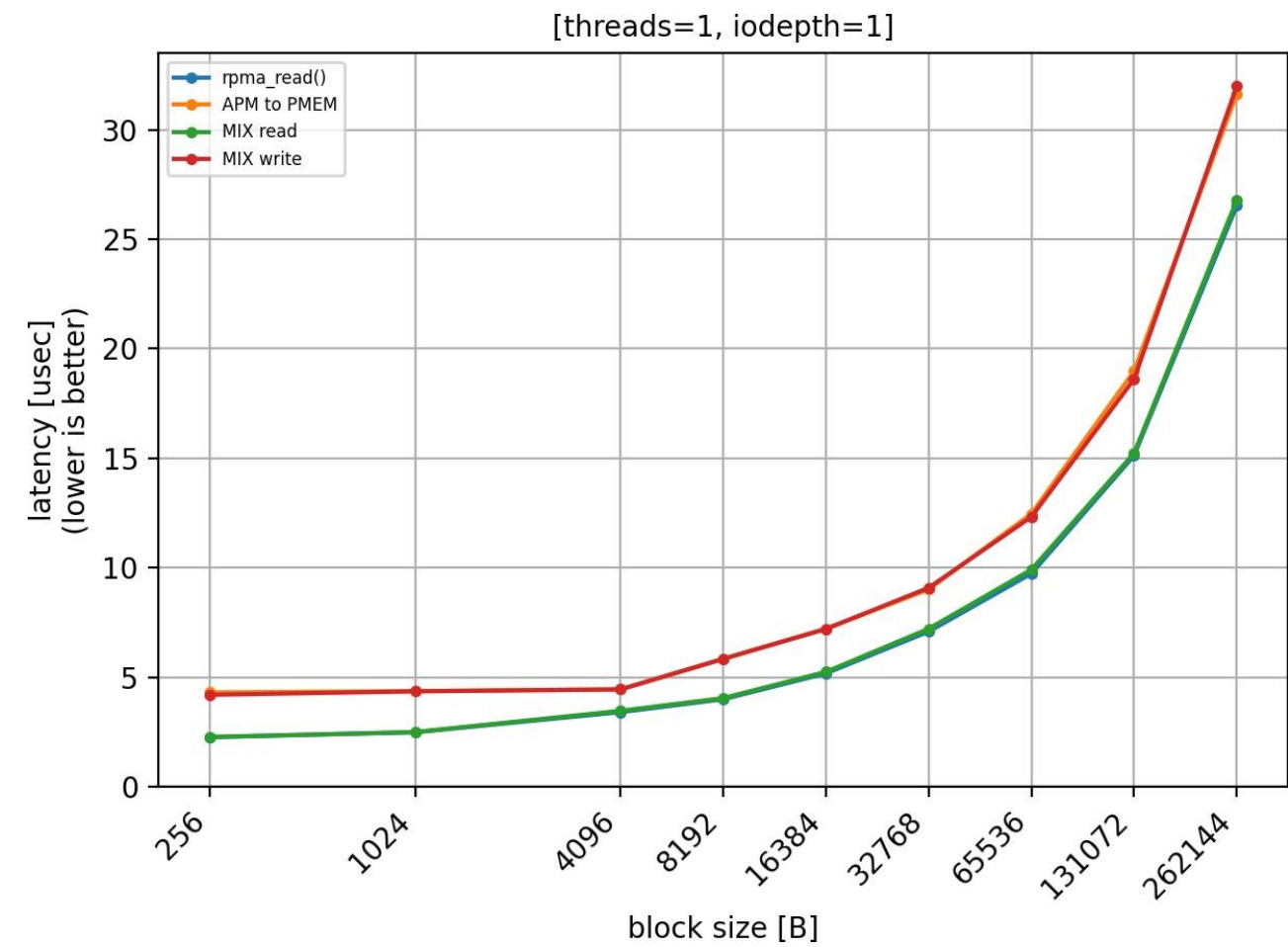
Benchmarking mixed reads from PMem (rpma\_read()) on the **RPMA Target** and writes (**APM**) to the same PMem on **the RPMA Target** side in the configured ratio (**MIX**) and comparing it to pure rpma\_read() (100% reads) and **APM** (100% writes).

## Mix against PMem: Latency

Comparing the latency of rpma\_read() and **APM** in the **MIX** workload vs the latency of 100% reads and 100% writes configurations.

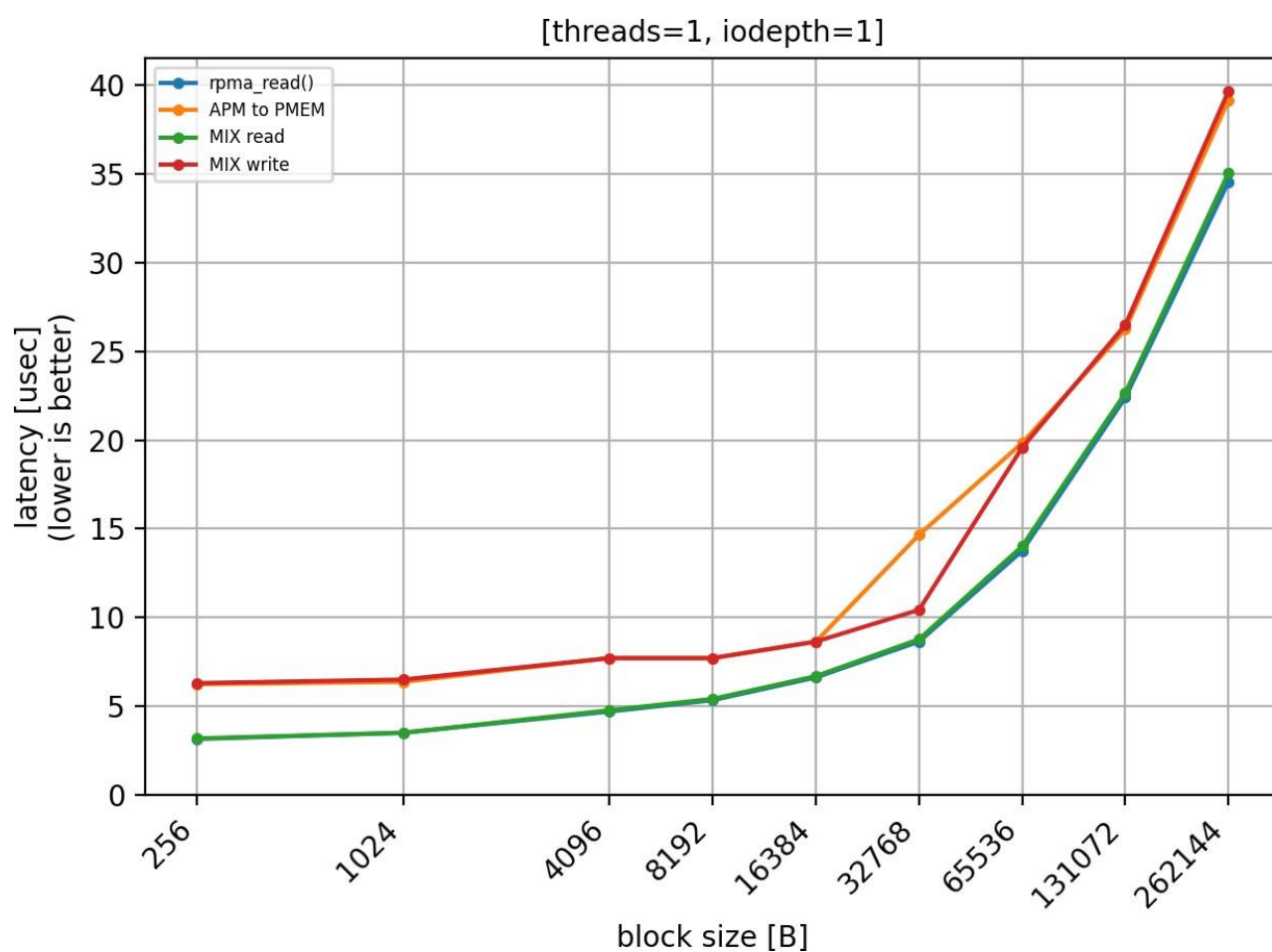
Item	Description
Server - FIO engine configuration	<pre>[global] ioengine=librpma_apm_server create_serialize=0 kb_base=1000 serverip=\$serverip port=7204 thread  [server] numjobs=1 direct_write_to_pmem=1 size=100MiB filename=/dev/dax/path</pre>
Client - FIO engine configuration	<pre>[global] ioengine=librpma_apm_client create_serialize=0 serverip=\$serverip port=7204 thread disable_clat=1 lat_percentiles=1 percentile_list=99.0:99.9:99.99:99.999  [client] sync=1 readwrite={read, write, rw, randread, randwrite, randrw} rwmixread=70 # valid only for the 'rw' and 'randrw' readwrite modes blocksize=\$blocksize ramp_time=15 time_based runtime=60</pre>

Figure 21. Latency (lat\_avg): MIX against PMEM vs rpma\_read() + APM to PMEM (seq)



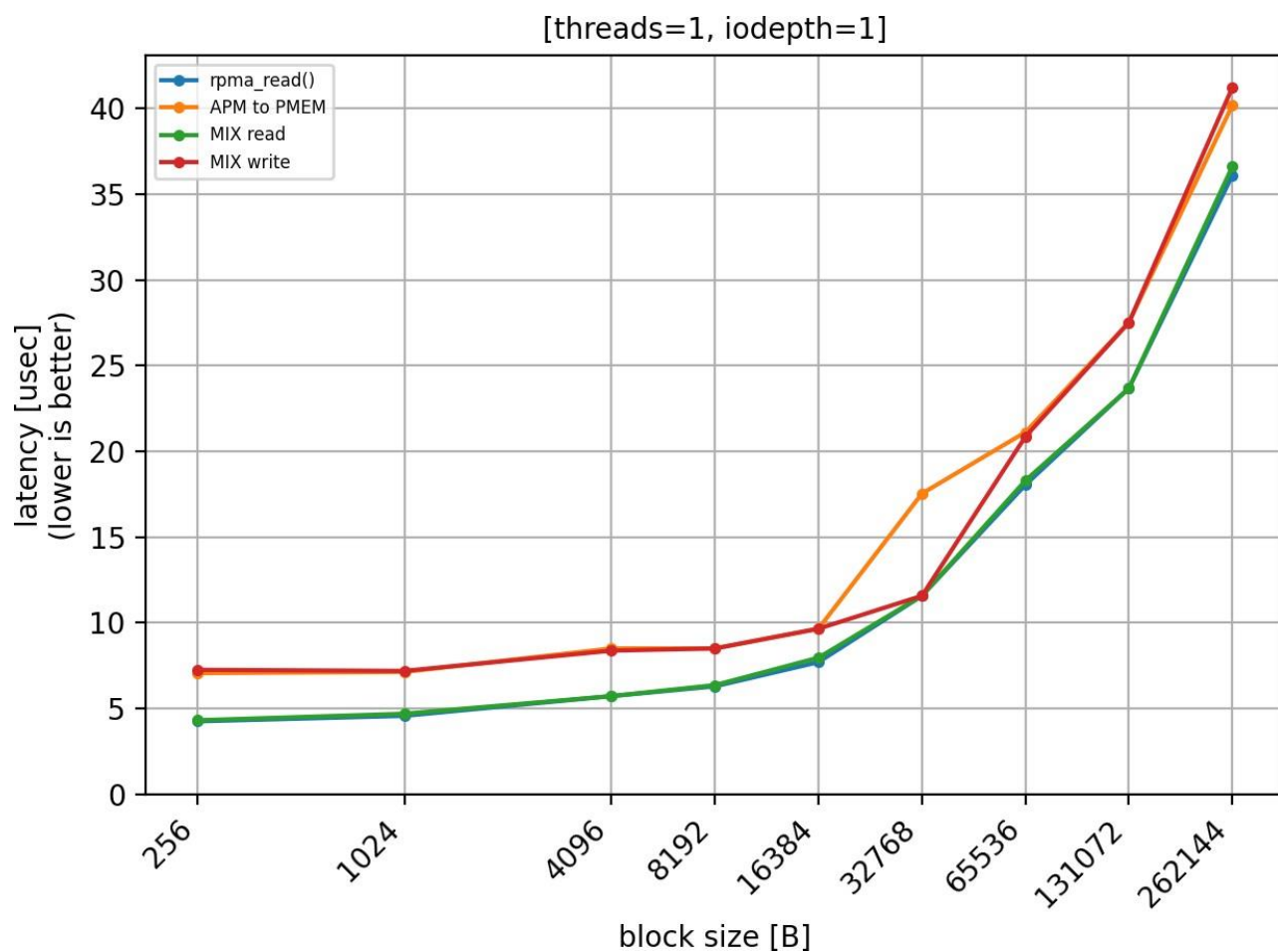
	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read()	2.26	2.48	3.40	3.99	5.17	7.08	9.73	15.12	26.54
APM to PMEM	4.30	4.36	4.44	5.84	7.21	9.01	12.46	18.98	31.65
MIX read	2.26	2.49	3.46	4.04	5.24	7.21	9.91	15.23	26.79
MIX write	4.19	4.35	4.44	5.83	7.20	9.08	12.32	18.61	32.00

Figure 22. Latency (lat\_pctl\_99.9): MIX against PMEM vs rpma\_read() + APM to PMEM (seq)



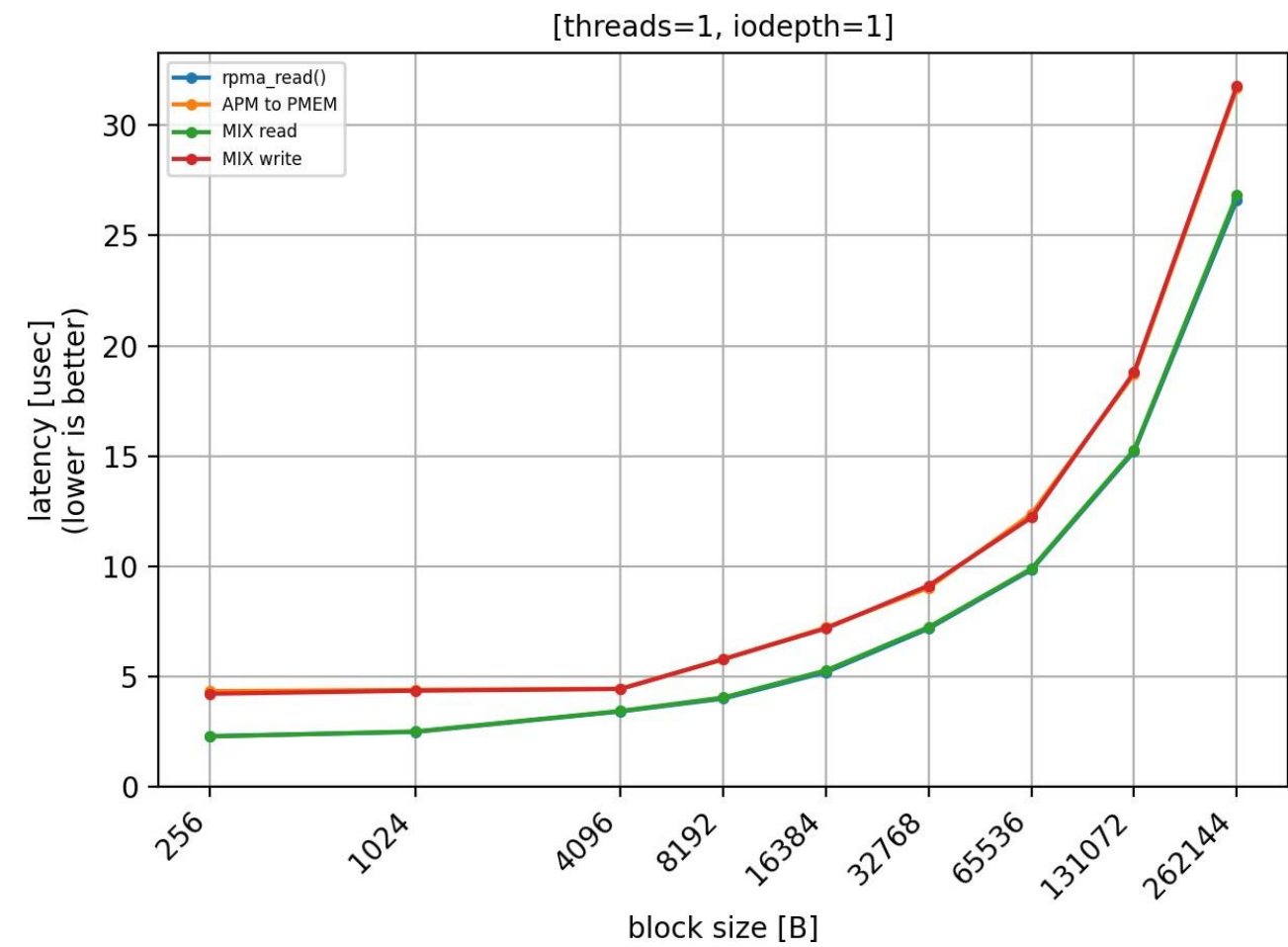
	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read()	3.15	3.50	4.70	5.34	6.62	8.64	13.76	22.40	34.56
APM to PMEM	6.24	6.37	7.71	7.71	8.64	14.66	19.84	26.24	39.17
MIX read	3.18	3.50	4.77	5.41	6.69	8.77	14.02	22.66	35.07
MIX write	6.30	6.50	7.71	7.71	8.64	10.43	19.58	26.50	39.68

Figure 23. Latency (lat\_pctl\_99.99): MIX against PMEM vs rpma\_read() + APM to PMEM (seq)



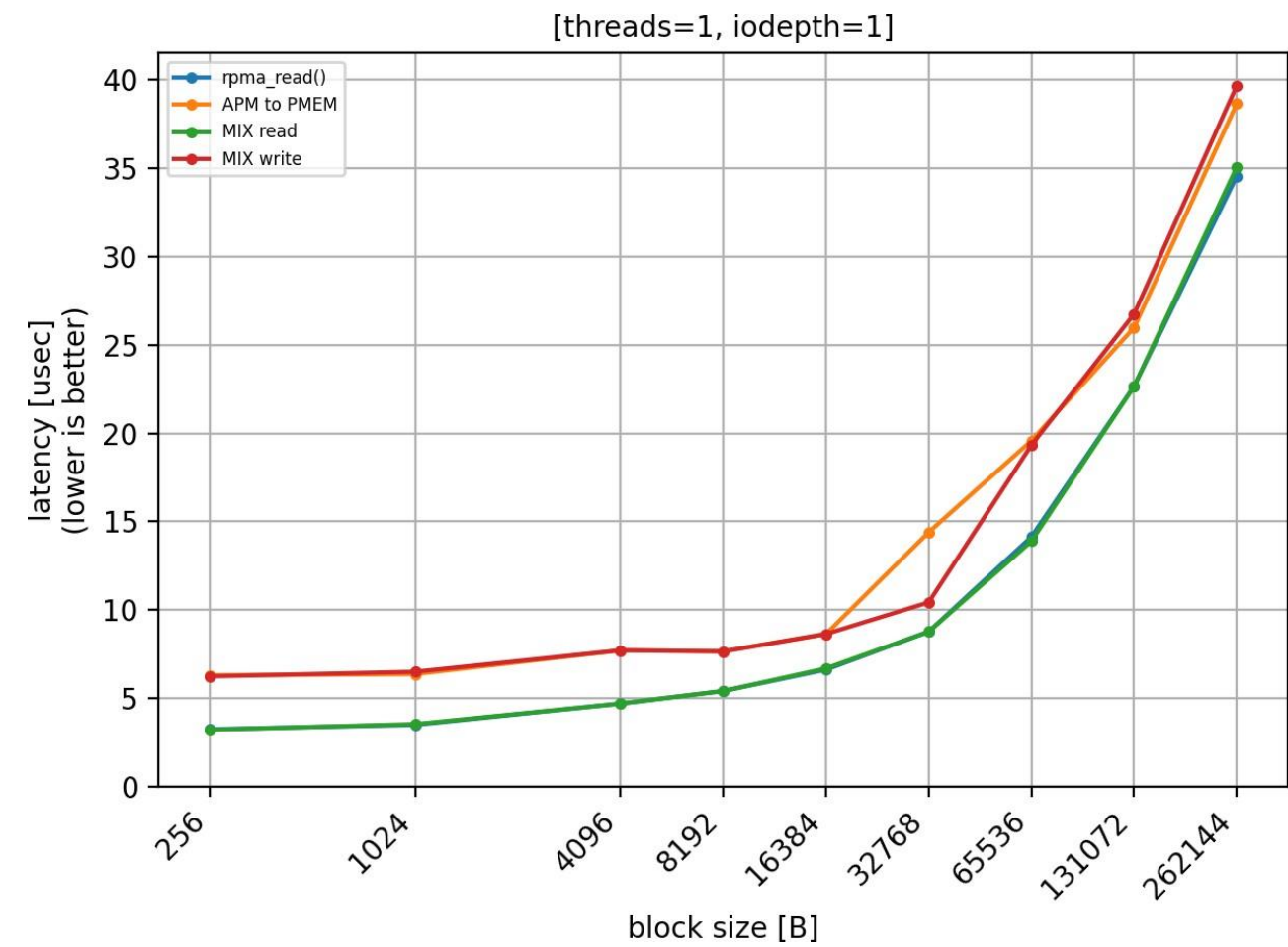
	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read()	4.26	4.58	5.73	6.30	7.71	11.58	18.05	23.68	36.10
APM to PMEM	7.07	7.14	8.51	8.51	9.66	17.54	21.12	27.52	40.19
MIX read	4.32	4.70	5.73	6.37	7.97	11.58	18.30	23.68	36.61
MIX write	7.26	7.20	8.38	8.51	9.66	11.58	20.86	27.52	41.22

Figure 24. Latency (lat\_avg): MIX against PMEM vs rpma\_read() + APM to PMEM (rand)



	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read()	2.29	2.48	3.41	3.99	5.19	7.17	9.83	15.20	26.62
APM to PMEM	4.33	4.39	4.44	5.78	7.22	9.02	12.37	18.70	31.68
MIX read	2.27	2.50	3.43	4.04	5.27	7.24	9.89	15.28	26.84
MIX write	4.21	4.35	4.43	5.78	7.18	9.12	12.21	18.80	31.79

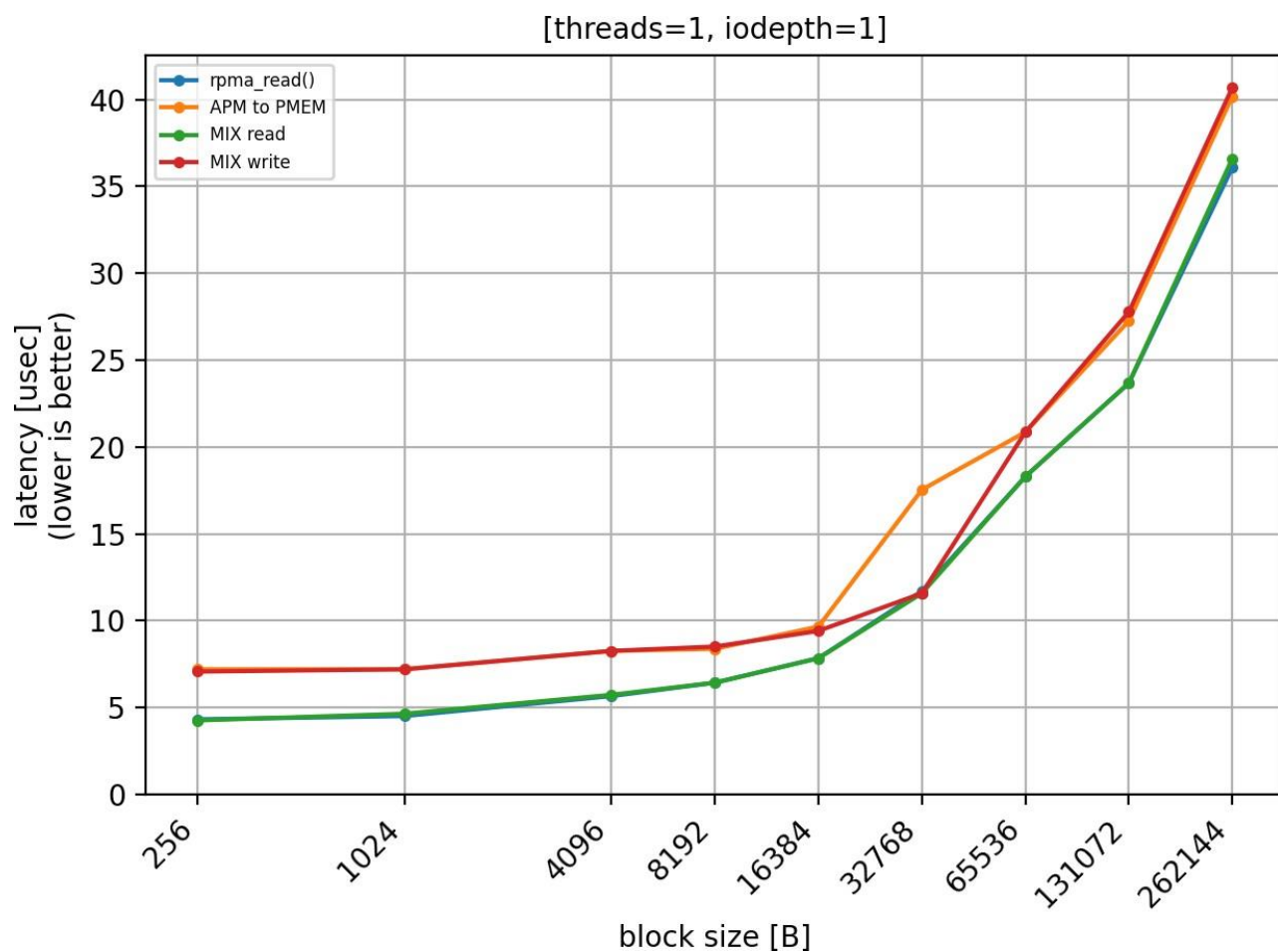
Figure 25. Latency (lat\_pctl\_99.9): MIX against PMEM vs rpma\_read() + APM to PMEM (rand)



	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read()	3.25	3.50	4.70	5.41	6.62	8.77	14.14	22.66	34.56
APM to PMEM	6.30	6.37	7.71	7.65	8.64	14.40	19.58	25.98	38.66
MIX read	3.22	3.54	4.70	5.41	6.69	8.77	13.89	22.66	35.07
MIX write	6.24	6.50	7.71	7.65	8.64	10.43	19.33	26.75	39.68



Figure 26. Latency (lat\_pctl\_99.99): MIX against PMEM vs rpma\_read() + APM to PMEM (rand)



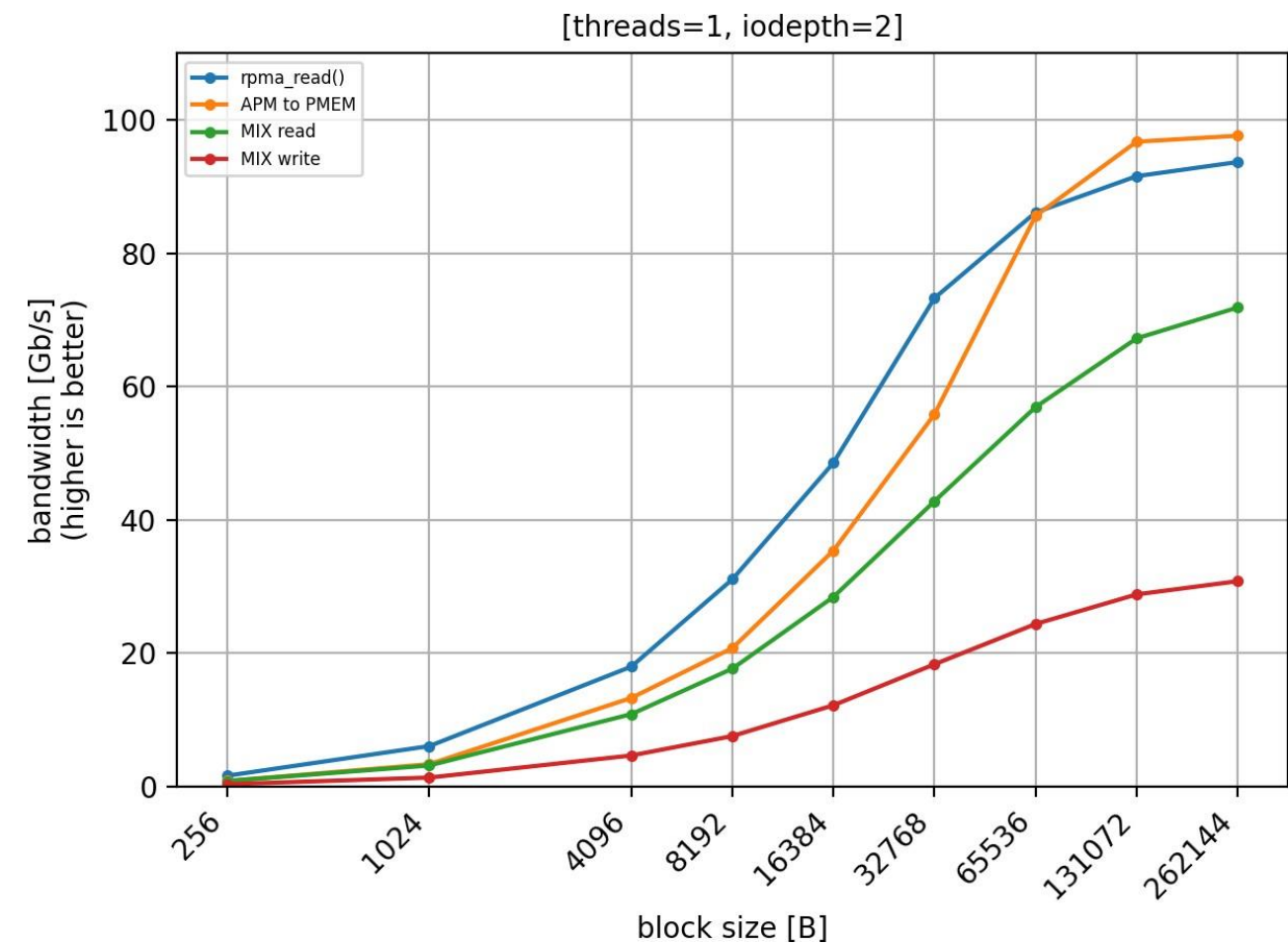
	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read()	4.32	4.51	5.66	6.43	7.84	11.71	18.30	23.68	36.10
APM to PMEM	7.20	7.20	8.26	8.38	9.66	17.54	20.86	27.26	40.19
MIX read	4.26	4.64	5.73	6.43	7.84	11.58	18.30	23.68	36.61
MIX write	7.07	7.20	8.26	8.51	9.41	11.58	20.86	27.78	40.70

## Mix against PMem: Bandwidth

Comparing the bandwidth of `rpma_read()` and **APM** in the **MIX** workload vs the latency of 100% reads and 100% writes configurations.

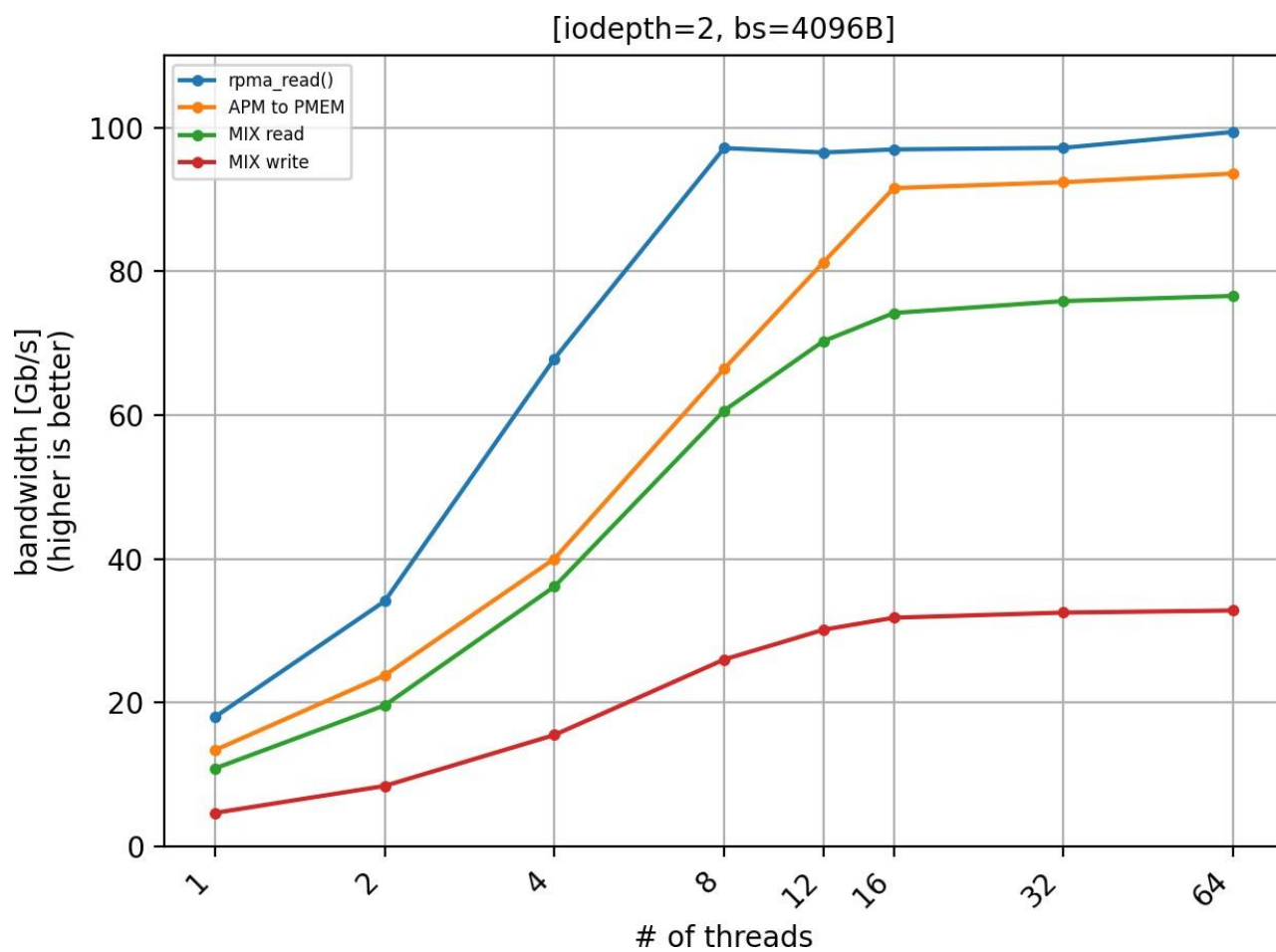
Item	Description
Server - FIO engine configuration	<pre>[global] ioengine=librpma_apm_server create_serialize=0 kb_base=1000 serverip=\$serverip port=7204 thread  [server] numjobs=\$numjobs direct_write_to_pmem=1 size=100MiB filename=/dev/dax/path</pre>
Client - FIO engine configuration	<pre>[global] ioengine=librpma_apm_client create_serialize=0 serverip=\$serverip port=7204 thread disable_clat=1 lat_percentiles=1 percentile_list=99.0:99.9:99.99:99.999  [client] numjobs=\$numjobs group_reporting=1 iodepth=2 readwrite={read, write, rw, randread, randwrite, randrw} rwmixread=70 # valid only for the 'rw' and 'randrw' readwrite modes blocksize=\$blocksize ramp_time=15 time_based runtime=60</pre>

Figure 27. Bandwidth (bs): MIX against PMEM vs rpma\_read() + APM to PMEM (seq)



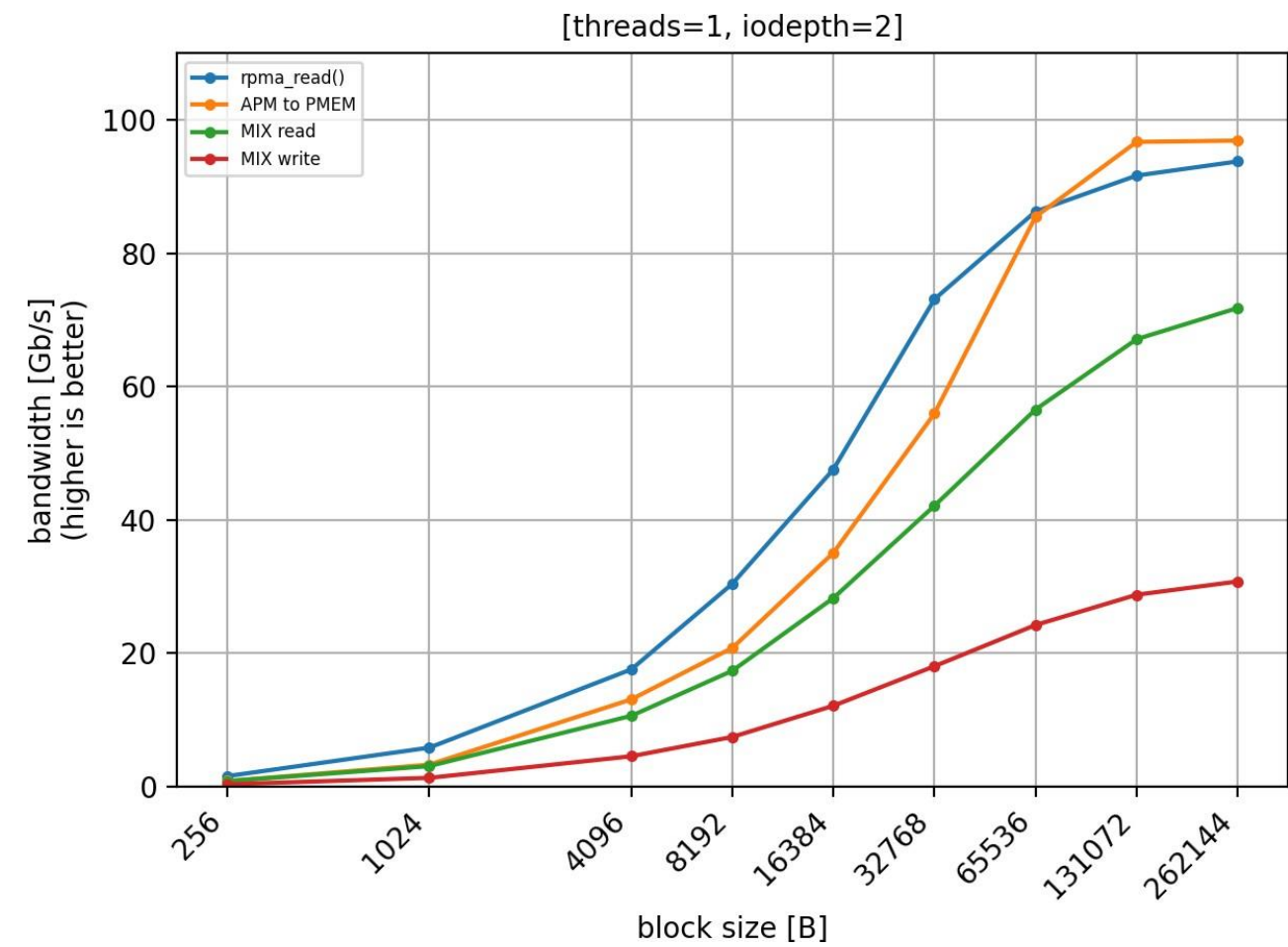
	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read()	1.63	6.07	18.00	31.10	48.58	73.29	86.13	91.58	93.70
APM to PMEM	0.85	3.36	13.30	20.78	35.41	55.86	85.62	96.75	97.65
MIX read	0.83	3.16	10.87	17.69	28.46	42.81	56.93	67.24	71.87
MIX write	0.36	1.35	4.66	7.58	12.20	18.35	24.39	28.83	30.81

Figure 28. Bandwidth (threads): MIX against PMEM vs rpma\_read() + APM to PMEM (seq)



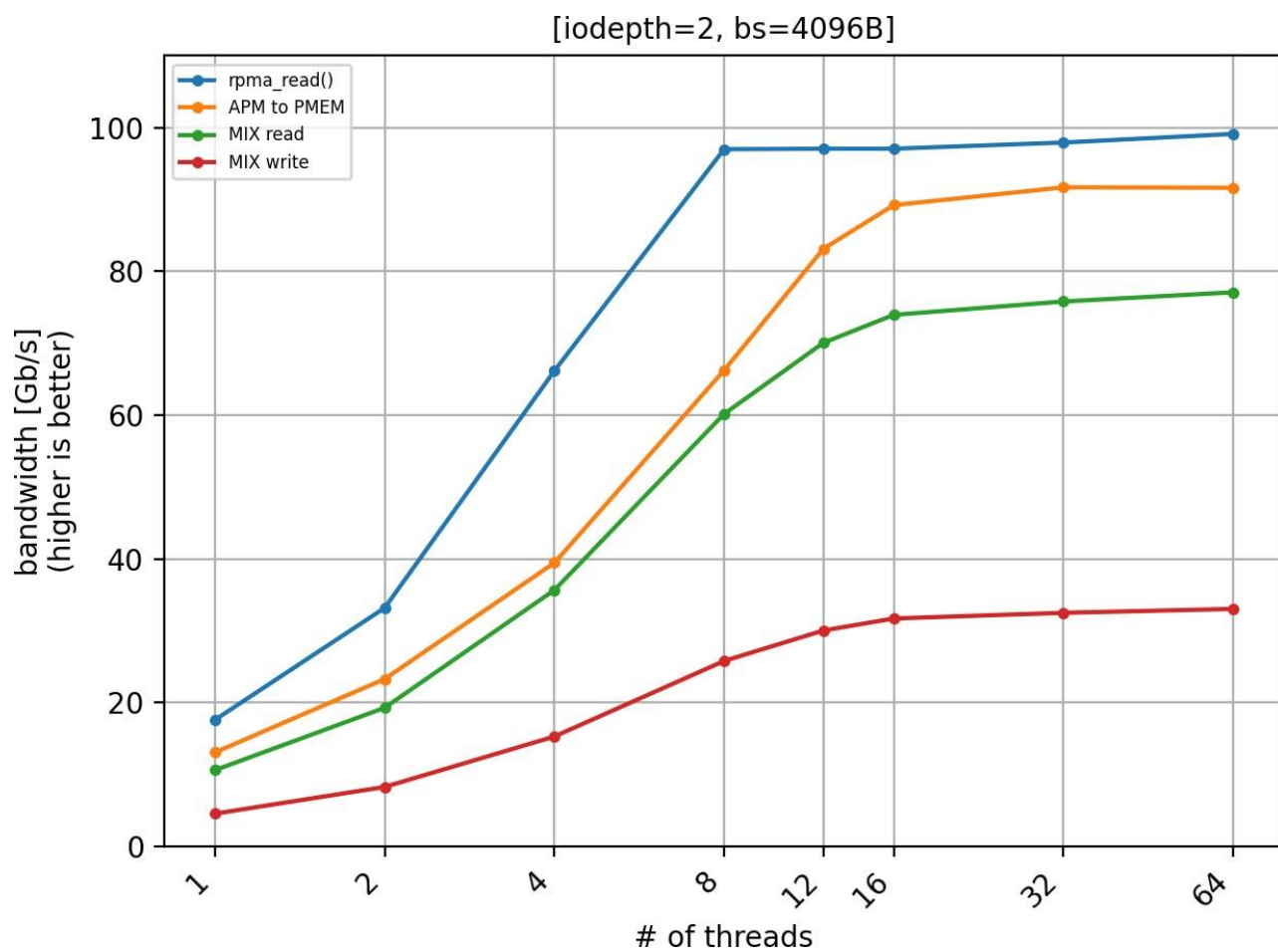
	1	2	4	8	12	16	32	64
rpma_read()	18.00	34.11	67.86	97.16	96.54	96.97	97.19	99.41
APM to PMEM	13.38	23.79	40.01	66.43	81.26	91.59	92.41	93.61
MIX read	10.84	19.61	36.14	60.63	70.29	74.21	75.87	76.58
MIX write	4.65	8.40	15.49	25.99	30.13	31.81	32.52	32.81

Figure 29. Bandwidth (bs): MIX against PMEM vs rpma\_read() + APM to PMEM (rand)



	256	1024	4096	8192	16384	32768	65536	131072	262144
rpma_read()	1.57	5.84	17.59	30.44	47.61	73.13	86.28	91.67	93.79
APM to PMEM	0.84	3.29	13.09	20.81	35.09	55.98	85.50	96.72	96.92
MIX read	0.81	3.07	10.62	17.38	28.30	42.14	56.57	67.12	71.80
MIX write	0.35	1.31	4.55	7.45	12.13	18.06	24.24	28.78	30.77

Figure 30. Bandwidth (threads): MIX against PMEM vs rpma\_read() + APM to PMEM (rand)



	1	2	4	8	12	16	32	64
rpma_read()	17.60	33.18	66.17	97.00	97.07	97.07	97.93	99.13
APM to PMEM	13.06	23.28	39.45	66.24	83.14	89.22	91.70	91.64
MIX read	10.60	19.28	35.66	60.16	70.06	73.95	75.81	77.09
MIX write	4.54	8.26	15.28	25.79	30.03	31.70	32.49	33.03

Ref: 2021\_08\_27\_CLX