# arm
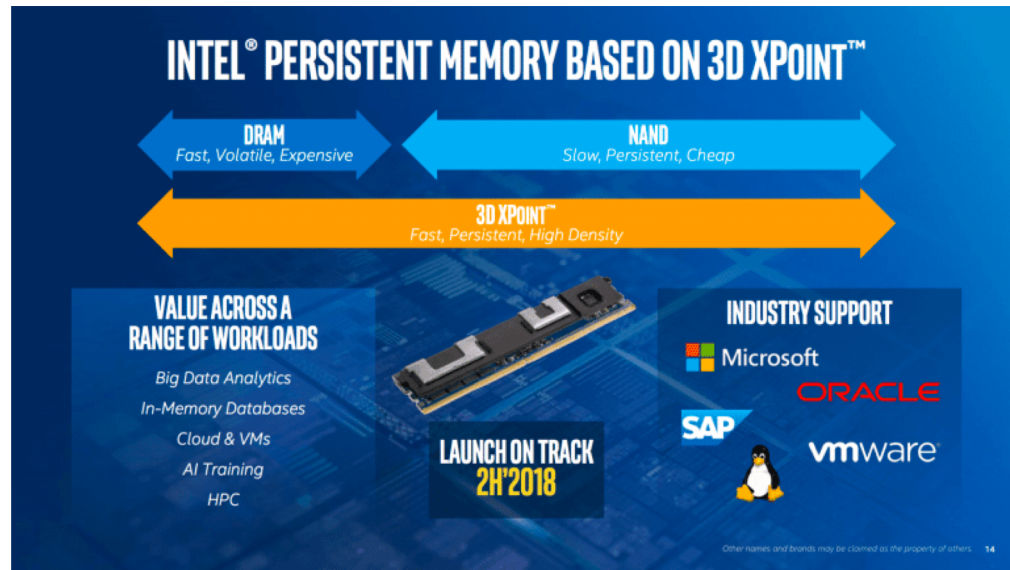
# *Quantifying the Performance Overheads of PMDK*

**William Wang**, Stephan Diestelhorst

2 October 2018

# NVDIMMs Shipping in 2H'2018



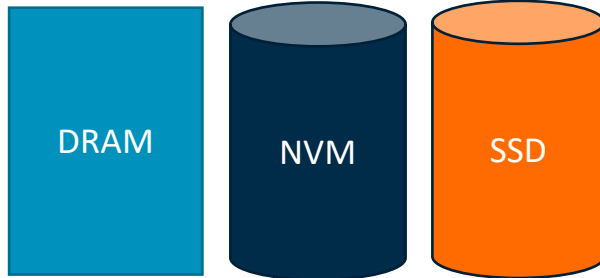INTEL® PERSISTENT MEMORY BASED ON 3D XPOINT™

DRAM
Fast, Volatile, Expensive

NAND
Slow, Persistent, Cheap

3D XPOINT™
Fast, Persistent, High Density

VALUE ACROSS A
RANGE OF WORKLOADS

Big Data Analytics
In-Memory Databases
Cloud & VMs
AI Training
HPC

LAUNCH ON TRACK
2H'2018

INDUSTRY SUPPORT
Microsoft
ORACLE
SAP
vmware

Other names and brands may be claimed as the property of others. 14

- 3DXP NVDIMM shipping 2H'2018
  - DDR4 -> $10/GB, NAND -> $1/4 per GB
  - 3DXP -> $2/GB, up to 512GB
- Arm NVDIMM-ready timeframe around 2022

arm

# Multiple System Use-Cases for NVM
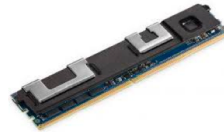
## Faster Storage
**1000x faster than NAND**



Storage
- Filesystem bottlenecks
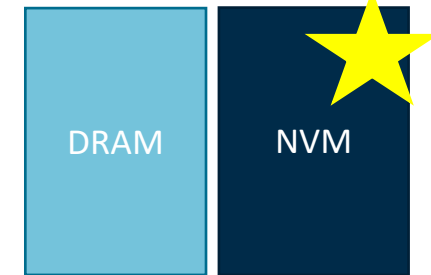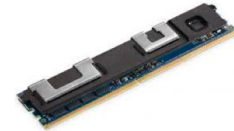
## Denser Mem
**10x denser than DRAM**



Transformative Capacity/TCO-advantage
- Endurance
- Bandwidth
- Caching

## Persistent Mem
**Non-Volatile**



Persistency
- Failure atomicity
- Ordering
- ISA

arm

# Spectrum of Application Performance Boost with NVM

## Executive summary

| Config→ | NVMe Flash SSD | Baseline<br>NVMe Flash SSD | Faster Storage<br>PMEM + Flash | Persistent Memory<br>PMEM-only |
|---|---|---|---|---|
| | AOF every *second*<br>+ RDB | AOF every *request*<br>+ RDB | AOF on PMEM<br>+ RDB on Flash | SW persistence<br>*no* RDB |
| **Data loss potential** | Up to 1 sec worth | None | None | None |
| **Crash restart** | Slowest | Slowest | Fast | Instantaneous |
| **Application rewrite** | No | No | No | Yes |
| **"Memory" cost^** | 1x | 1x | 1.25x§ | 0.27x† |
| **Iso-SLA performance(‡)** | 1x | 0.165x (0.149x) | 1.46x (1.31x) | 2.18x (1.96x) |

**13x**

- **13x** *perf boost,* **1/4** *mem cost*

- *No data loss*

- *Instantaneous restart*

- ***Rewrite application***

^ Flash costs are marginal for all configurations, thus ignored.
§ Assumes 8GB/core DRAM + 1 GB/core SCM capacity – SCM capacity very generous for log use-case, previous studies used as low as 40MB/core.
† Assumes 8GB/core DRAM capacity, SCM = 0.25x DRAM $/GB, 1:8 capacity ratio for DRAM:SCM cache ratio.
‡ Measured performance de-rated by 0.9x to account for a slower-than-DRAM SCM media fronted by a DRAM cache.
* Unable to achieve same SLA as baseline with p99 for the every-request AOF logging configuration – hence the unconstrained SLA comparison.
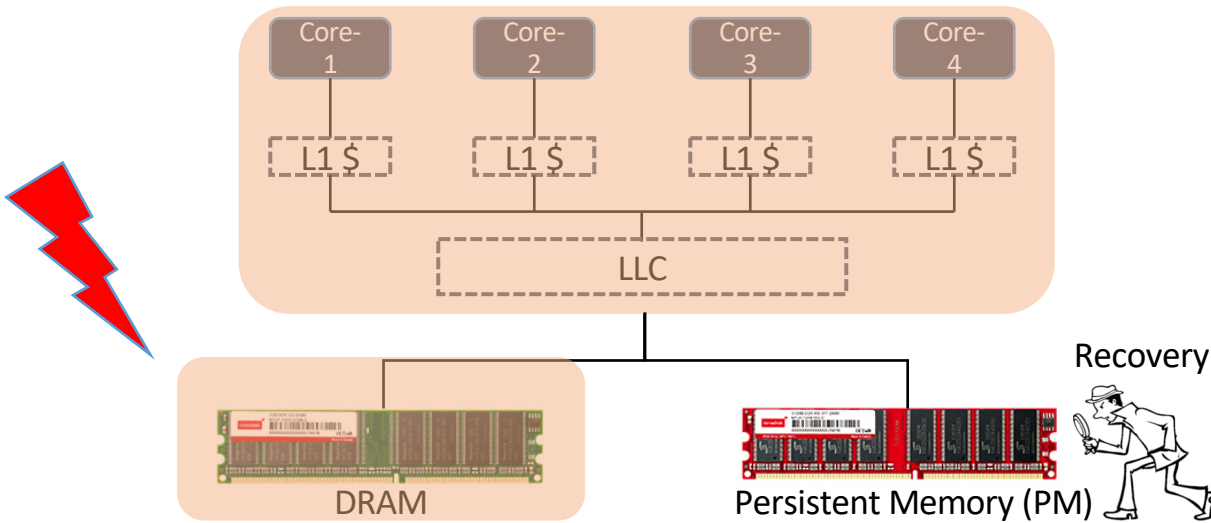
arm

Slide Courtesy of Kshitij Sudan

arm

# Why Application Rewrite



**Recovery can inspect the data-structures in PM to restore system to a consistent state**

- Crash consistency (failure atomicity) is needed to ensure recovery can restore system to a consistent state

  - Data move through volatile memories before they get written to PM

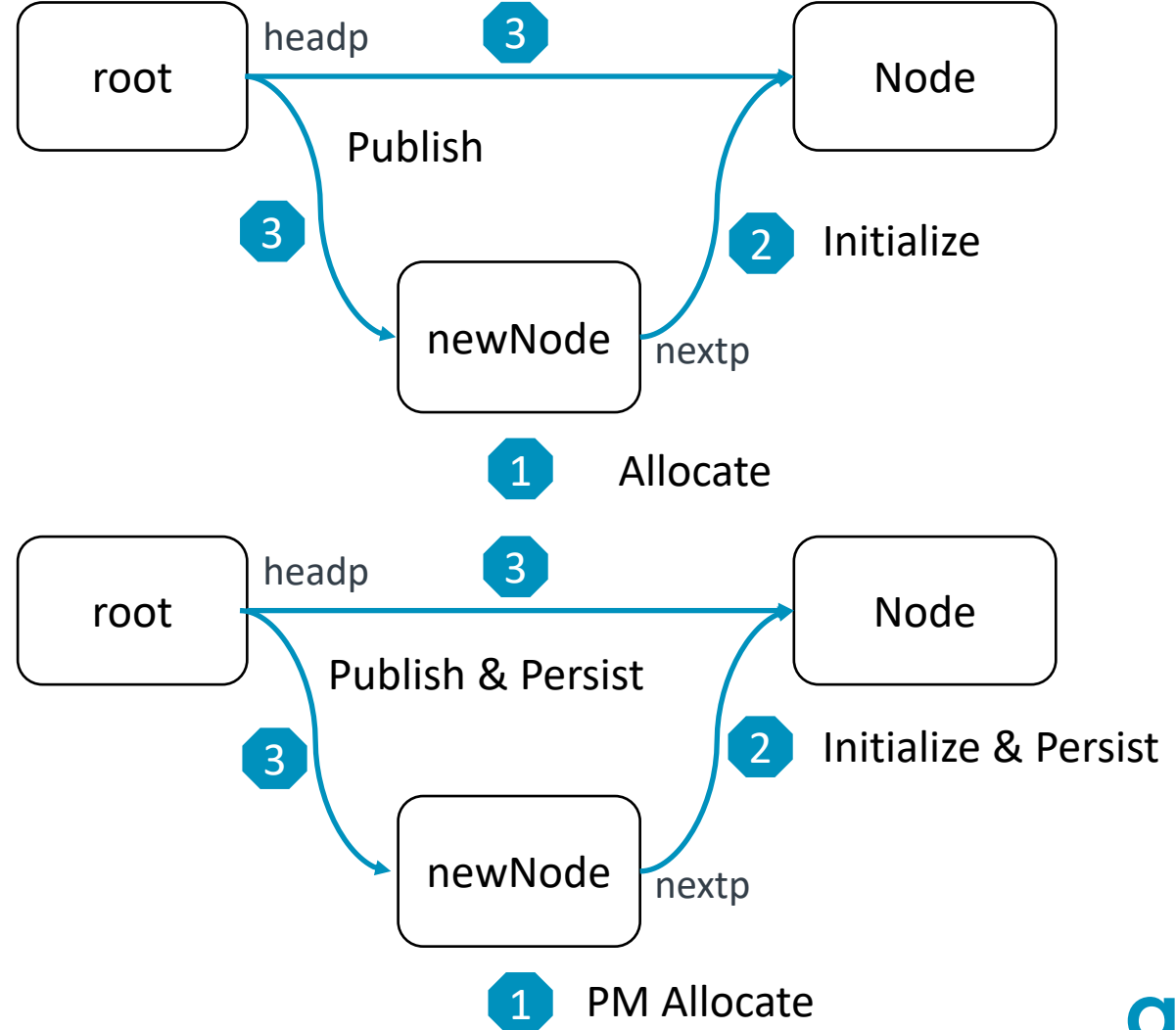  - Using CPU cache flushes and fence instructions

arm

# Example: Add a Node to a Linked List

```
1   // Add a node to a linked list
2   void
3   addnode(struct root *rootp, int data)
4   {
5       struct node *newnodep;
6       if ((newnodep = calloc(1,
7           sizeof(struct node))) == NULL)
8         fatal("out of memory");
9       newnodep->data = data;
10      newnodep->nextp = rootp->headp;
11      rootp->headp = newnodep;
12  }
```

```
1   // Add failure atomicity
2   void
3   addnode(struct root *rootp, int data)
4   {
5       struct node *newnodep;
6       int flag = 0;
7       if ((newnodep = pm_calloc(1,
8           sizeof(struct node))) == NULL)
9         fatal("out of memory");
10      newnodep->data = data;
11      newnodep->nextp = rootp->headp;
12      pm_flush(newnodep,
13          sizeof(struct node));
14      pm_fence();
15      rootp->headp = newnodep;
16      pm_flush(&(rootp->headp),
17          sizeof(rootp->headp));
18      pm_fence();
19  }
```

# Example: Add Concurrency to the Mix

```
 1   // Add a node to a linked list
 2   void
 3   addnode(struct root *rootp, int data)
 4   {
 5       struct node *newnodep;
 6       if ((newnodep = calloc(1,
 7           sizeof(struct node))) == NULL)
 8           fatal("out of memory");
 9       newnodep->data = data;
10       newnodep->nextp = rootp->headp;
11       rootp->headp = newnodep;
12   }
```

```
 1   // Add multithread atomicity
 2   void
 3   addnode(struct root *rootp, int data)
 4   {
 5       struct node *newnodep;
 6       if ((newnodep = calloc(1,
 7           sizeof(struct node))) == NULL)
 8           fatal("out of memory");
 9       newnodep->data = data;
10       /* lock the critical section */
11       pthread_mutex_lock(
12           &rootp->listlock);
13       newnode->nextp = rootp->headp;
14       rootp->headp = newnodep;
15       pthread_mutex_unlock (
16           &rootp->listlock);
17   }
```

```
 1   // Add failure atomicity
 2   void
 3   addnode(struct root *rootp, int data)
 4   {
 5       struct node *newnodep;
 6       if ((newnodep = pm_calloc(1,
 7           sizeof(struct node))) == NULL)
 8           fatal("out of memory");
 9       newnodep->data = data;
10       /* lock the critical section */
11       pthread_mutex_lock(
12           &rootp->listlock);
13       newnodep->nextp = rootp->headp;
14       pm_flush(newnode,
15           sizeof(struct node));
16       rootp->headp = newnodep;
17       pm_flush(&(rootp->headp),
18           sizeof(rootp->headp));
19       pthread_mutex_unlock (
20           &rootp->listlock);
21   }
```
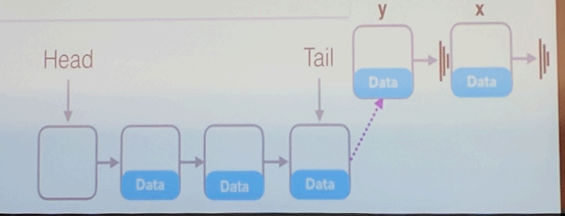
*What can still go wrong?*

Fences are missing after persists, they can get persisted out of order!!!

*What about making this lock-free?*



## It's Complicated

- *Oracle Labs at PPoPP'18*

Source: M. Friedman, A Persistent Lock-Free Queue for Non-Volatile Memory, PPoPP'18

   Simplify lock-based persistent programming "**Persistency for Synchronization-free Regions**" in **PLDI '18**

arm

# Persistent Memory Programming Models

**Native Persistence**

```
pt->x = 1;
pt->y = 1;
dccvap(&pt->x)
dccvap(&pt->y)
dsb

flag=1;
dccvap(&flag)
dsb
```

**Library Persistence – Atomic**

```
pt->x = 1;
pt->y = 1;
pmem_persist(&pt,
sizeof(pt))

flag = 1;
pmem_persist(&flag,
sizeof(flag))
```

**Library Persistence – Durable TXs**

```
TX_BEGIN{
pt->x = 1;
pt->y = 1;
} TX_END
```

Programming simpler, overhead higher

**arm**

# Add Concurrency to the Mix

## Lib Persistence – Lock-free

```
1   // ctor, dtor, copy and assign
2   struct Node
3   {
4       int x;
5       int y;
6   };
7
8   //lock-free update
9   int updateNode(Node* pt, int flag)
10  {
11      // allocate a new node (CoW)
12      atomic<Node*> newpt = new Node();
13      persist(newpt);
14
15      // update the node
16      newpt->x = 1;
17      newpt->y = 1;
18      persist(&newpt->x);
19      persist(&newpt->y);
20
21      // swing the pointer to the new Node and update flag
22      while(true) {
23          // copy on write
24          atomic<Node*> cowpt = pt;
25
26          if(CAS(&pt, cowpt, newpt)) {
27              persist(&pt);
28              // no contention for flag
29              flag = 1;
30              persist(&flag);
31              return 1;
32          }
33      }
34  }
```

## Lib Persistence – Lock-based

```
mutex.lock()
pt->x = 1;
pt->y = 1;
pmem_persist(&pt,
sizeof(pt))

flag = 1;
pmem_persist(&flag,
sizeof(flag))
mutex.unlock()
```
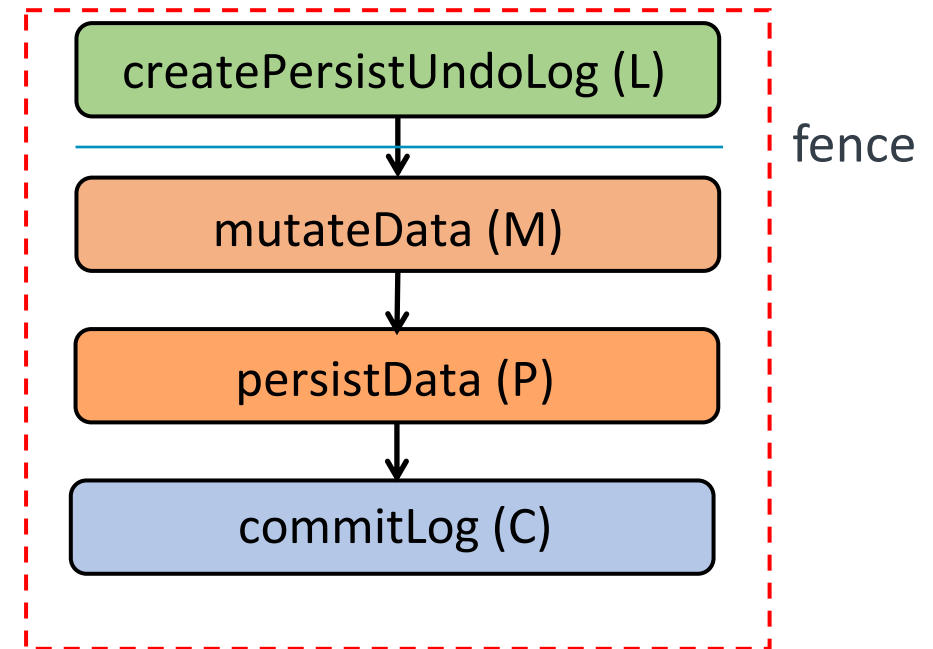
## Lib Persistence – Durable TXs

```
TX_BEGIN{
pt->x = 1;
pt->y = 1;
} TX_END
```

Programming simpler, overheads higher

arm

# Durable Transactions

- TXs provide clean failure semantics

- Accelerate w. Persistent HTM

```
TX_BEGIN{
pt->x = 1;
pt->y = 1;
} TX_END
```
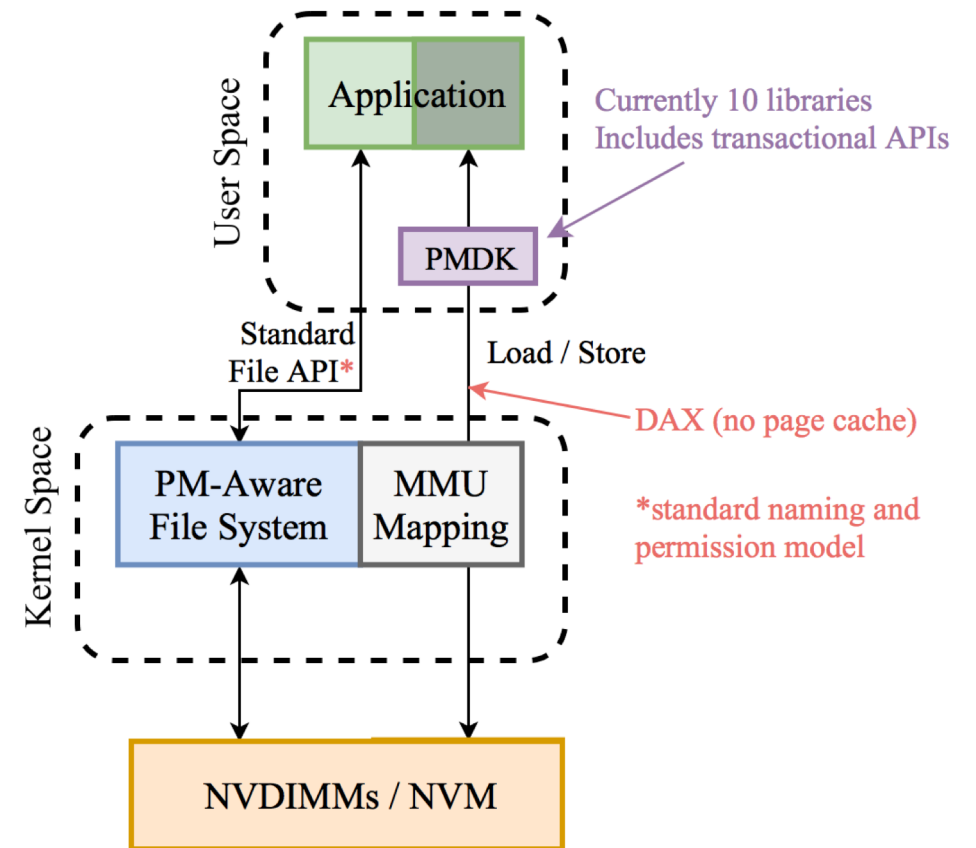
- Durable transactions guarantee Tx failure atomicity

  - Ordered persists + logging

# PMDK (Persistent Memory Development Kit)

Formally NVML, 'pmem libraries'

- PMDK provides transactional APIs for persistent memory programming

  - libpmemobj transactional APIs

  - Use fine-grained logging and cache flushes

- Works on 64-bit Linux, Windows and 64-bit FreeBSD

- Supports x86 and Armv8 (experimental)



Ref: pmem.io

arm

# Analyze Cost of Durable Transactions

## Logging, Persisting and Ordering

**ARM**

# Cost of Durable Transactions

Logging, Persisting and Ordering

## Persisting

- Latency too high if PoP is off-chip

- Undo log needs to be persisted in TX

- Instruction bloat due to looping through addresses with DCCVAP

## Logging

- Write amplification, i.e., write twice

- Copying (costly, and cache pollution)

- Barriers with undo logging

- Memory fragmentation

- Code bloat due to logging

## Ordering

- Serialize memory accesses

- Excessive # of barriers via API calls

- Persisting barriers same as ordering barriers

arm

# Performance Measurement Setup
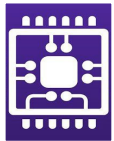
**Payload**

Redis-Benchmark

TPCC

**Workload**

Redis

Maps (PMDK/NVML)

**CPU**

CLFLUSH vs. CLFLUSHOPT

Fixed frequency and ASLR off
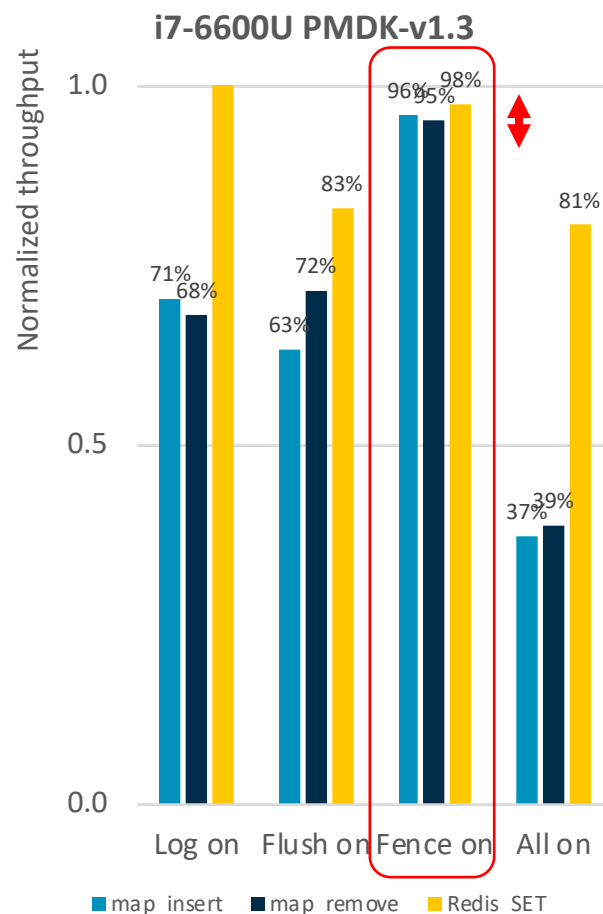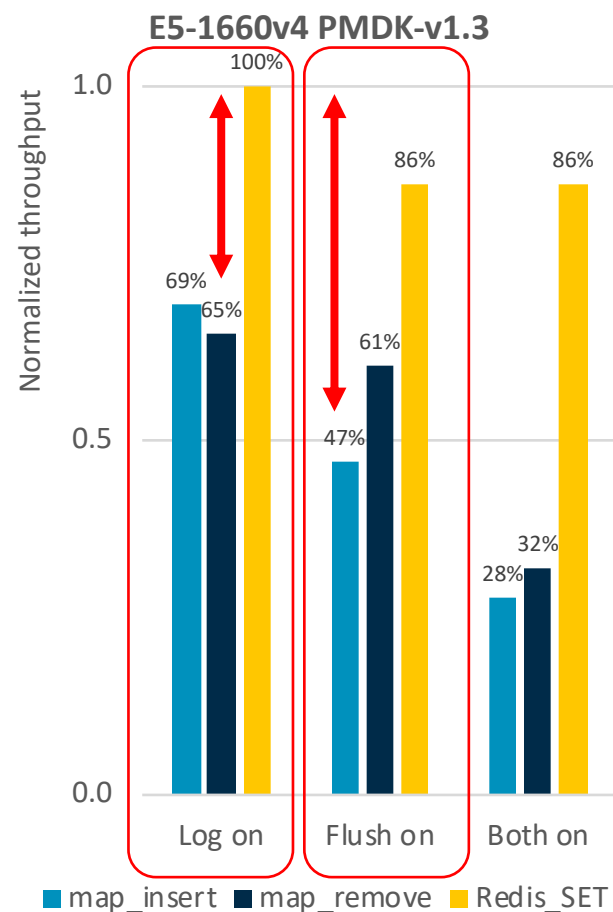Thread and memory affinity in multi-node

**Memory**

Local DRAM

Remote DRAM to emulate NVM
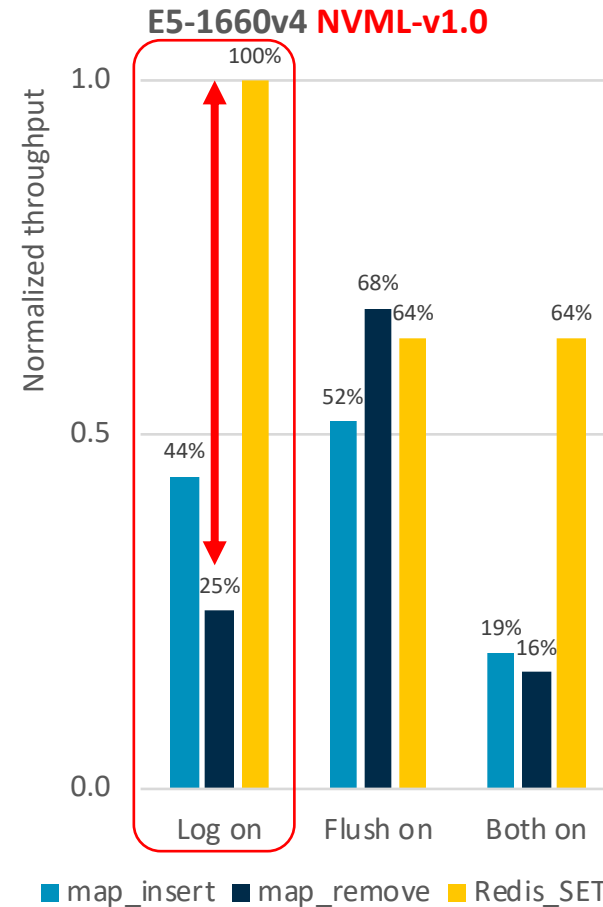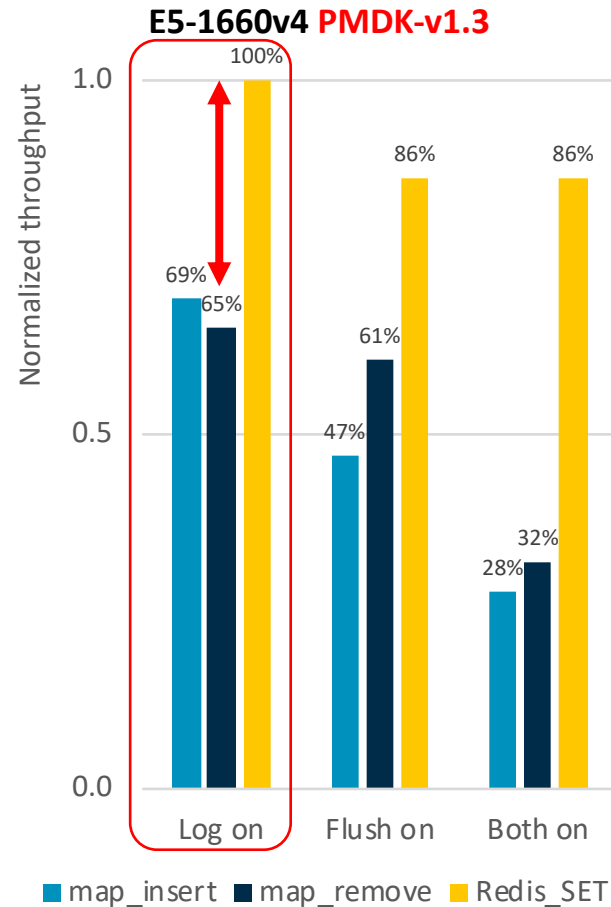
arm

# Flushing, Logging and Fencing Overhead



- Logging overhead 0~35%
- Flushing overhead 14~53%
- Fencing overhead 2~5%

Baseline: PMDK without flushing/fencing and logging on

- Workloads: Map insert/remove, Redis Set. Implemented with NVML v.10 and v1.3 libpmemobj transactions
- Platforms: Intel E5-1660v4 with CLFLUSH and Intel i7-6600U with CLFLUSHOPT, single node with local DRAM

# Software Opt Can Reduce Logging Overhead

**E5-1660v4 PMDK-v1.3**

**E5-1660v4 NVML-v1.0**

- Software opt can reduce logging overhead

  - NVML v1.3 (left) reduces logging overhead over NVML v1.0 (right)

  - Redis-SET performance not affected by logging overhead

Baseline: PMDK without flushing/fencing and logging on

© 2018 Arm Limited    • Workloads: Map insert/remove, Redis Set. Implemented with NVML v.10 and v1.3 libpmemobj transactions
                         • Platforms: Intel E5-1660v4 with CLFLUSH and Intel i7-6600U with CLFLUSHOPT, single node with local DRAM

arm

# Summary of Performance Overheads Analysis

- Flushing overhead up to 53%.

- Logging overhead up to 35%.

- Fencing overhead up to 5%

- These overheads can be reduced by

  - Software optimizations

  - Hardware accelerations

**arm**

# Ease-of-Programming vs Performance

● : concurrent

■ : concurrent & durable

Performance

Lock-free

CAS ●
→
■ PCAS

HTM
●
→
■
PHTM

Fine-grained Locking ●
→
SFR ■

● STM

Ease-of-programming

| Programming Model | Ordering | Logging | Persisting | Translation |
|---|---|---|---|---|
| Atomics | ✓ | X | ✓ | ✓ |
| TM | ✓ | ✓ | ✓ | ✓ |
| Locking | ✓ | ✓ (log elision) | ✓ | ✓ |

## Ease-of-Programming

Persistent CAS as drop-in replacement for CAS

Persistent HTM as drop-in replacement for HTM

Compiler instrumented failure atomicity for locking

## Performance

Relaxed memory persistency model

Hardware accelerations

Software optimizations

arm

Thank You!

Danke!

Merci!

谢谢!

ありがとう!

Gracias!

Kiitos!

감사합니다

धन्यवाद

arm