# Operating System Support for Multi-Tiered Memory
## (A Case Study of Intel's Optane DCPMM)

Jeongseob Ahn

E-mail: jsahn@ajou.ac.kr
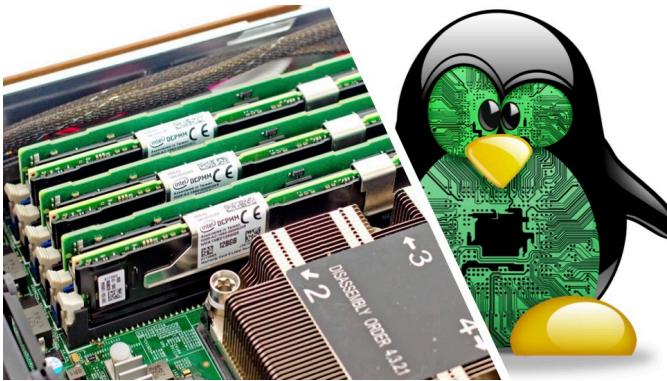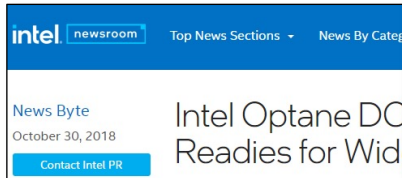Web: https://jeongseob.github.io

AJOU UNIVERSITY

# In this talk…

- Exploring the Design Space of Page Management for Multi-Tiered Memory Systems
  - Jonghyeon Kim, Wonkyo Choe, and Jeongseob Ahn
  - USENIX Annual Technical Conference (ATC), July 2021

- A Study of Memory Placement on Hardware-assisted Tiered Memory Systems
  - Wonkyo Choe, Jonghyeon Kim, and Jeongseob Ahn
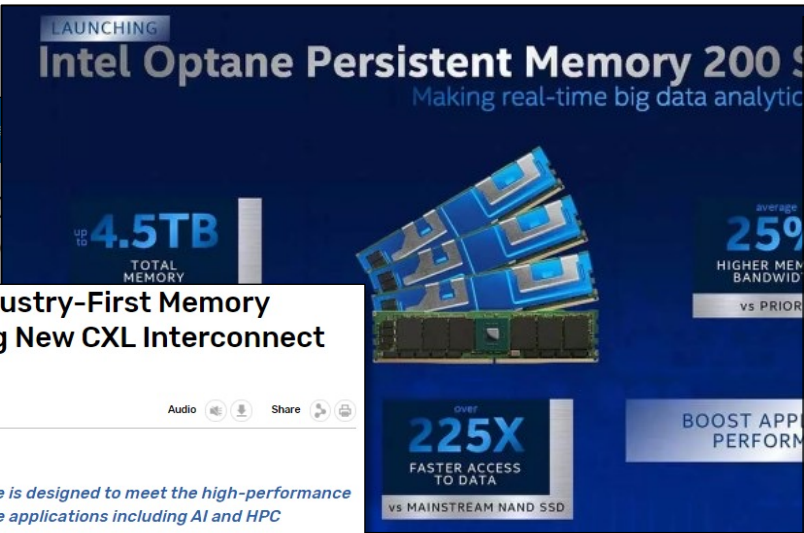  - IEEE Computer Architecture Letters (CAL), 19(2), July-December 2020

Artifact available at https://github.com/csl-ajou/AutoTiering

# Large-scale memory systems

# A multi-tiered memory system



- Modern server systems have formed NUMA
- DRAM and DCPMM share the memory controller
- Since Linux 5.0, Optane DCPMM can be exposed as a normal RAM
  - A DCPMM memory node is treated as a CPU-less NUMA node

(a)

(b)

SW-managed tiered memory organization
(A DCPMM memory node is treated as a CPU-less NUMA node)

# Memory access latency and bandwidth



4-socket `DRAM-based` NUMA machine

Multi-tiered NUMA machine

The critical factor in performance is not only **access locality** but also **access tier** of memory

# Default memory placement: `local-first`

State-of-the-art Linux kernel* has not considered the characteristics of fast (DRAM) and slow (DCPMM) memory with NUMA properties



Q. Why don't you reorder the fallback node list according to the actual performance?

# Limited page placement in current Linux
## Why not `AutoNUMA`?



Page movement to `CPU-less` nodes (DCPMM) is prohibited in the current Linux

Page movement is allowed only <u>when the target node has a free space</u> in the current Linux

# Need for page placement for tiered memory



graph500

Stock Linux 5.3 version

# Problems with current page management

## Problems (or limitations)

1. Allocation fallback does not consider access tier.

2. Pages are not promoted when upper-tier is full

3. Pages are never demoted or reclaimed to lower-tier memory

4. Page classification is too coarse-grained (binary: active or inactive)

# Exploiting access tier first and then locality

Conservative Promotion or Migration → `AutoTiering-CPM`



Case-1: page promotion

Case-2: page promotion or migration

`AutoTiering-CPM` provides alternatives for page migration failure due to fully occupied target memory node, leading to performance improvement

However, the upper-tier (DRAM) memory can still hold infrequently accessed data while frequently used pages reside in the lower-tier (DCPMM) memory

# Enforcing page promotion and migration

Opportunistic Promotion or Migration → `AutoTiering-OPM`



Fault page  Least accessed page

- Finding the least accessed page (LAP)
  1. Inactive page from `file-backed` region
  2. LAP page from `anonymous` region



*Access history is collected from AutoNUMA framework*

With `AutoTiering-OPM`, we can achieve better utilization of the upper-tier memory

# Hiding latency of page eviction
## A software optimization comes to rescue



Foreground / Background

Threads — CPU-0 (Memory Controller) — Interconnection — CPU-1 (Memory Controller)

DRAM (node-0) *fully occupied*

DRAM (Node-1)

DCPMM (node-2)

DCPMM (node-3)

Fault page / Least accessed page

**2b** Page eviction — DRAM read, DCPMM write → **3** Page promotion — DCPMM read, DRAM write

Reserved page / Least accessed page

**2** Page promotion — DCPMM read, DRAM write → **Critical path is shortened**

**2** Page eviction — DRAM read, DCPMM write → **3** Page return

# Experimental environments

- System
  - Intel(R) Xeon Gold 5218 CPU @ 2.30GHz x 2
  - 16GB DRAM x 2
  - 128GB Intel Optane DCPMM x 2
  - Linux kernel 5.3 with Ubuntu 18.04

- Benchmarks
  - SPECAccel (OpenMP)
  - GraphMat (PageRank)
  - Graph500 (BFS)
  - Liblinear

# Performance evaluation



BD: background demotion

Legend: Baseline (Stock Linux 5.3) | CPM | OPM (BD)

graph500: 6.99x
GraphMat: 2.31x
Liblinear: 3.69x
503.postencil: 2.48x
553.pclvrleaf: 2.42x
559.pmniGhost: 1.69x, 2.61x
560.pilbdc: 2.27x

* Most benchmarks are improved by `AutoTiering`

* With `CPM`, speedup is up to 2.48x in `503.postencil`

* With `OPM(BD)`, speedup is up to 6.99x in `graph500`

# Effectiveness of `AutoTiering-CPM`



Distribution of Memory Usages

`AutoTiering-CPM` makes better use of multi-tiered memory

# Effectiveness of `LAP` classification



Less effectively utilized

More effectively utilized

Stock Linux 5.3 version

AutoTiering-OPM

AutoTiering-OPM can promote frequently accessed pages
while demoting least accessed pages

# Effect of hiding demotion latency

Measured page promotion latency (CDF) with `ftrace`
- `OPMX`: Opportunistic Page Migration with Exchange*



We can reduce the promotion latency by deferring the page demotion as background

*Nimble Page Management for Tiered Memory Systems [ASPLOS '19]

# Is page exchange acutally needed?



Page exchange scheme*

instead of copying data into new pages, we transfer data between each pair of pages using copy thread(s) that use CPU registers as the temporary storage for in-flight iterative data exchange operations. This use of registers allows our mechanism to avoid allocating a complete temporary page.

* Nimble Page Management for Tiered Memory Systems [ASPLOS '19]

# Performance comparison with prior work

# Summary

- Commodity OSes are not mature enough to support multi-tiered memory systems

- We explored new page placement schemes to extract the full benefits of multi-tiered memory systems

- Future work
  - Redesigning the kernel thread demoting page migration to DCPMM with the consideration of limited memory bandwidth of DCPMM
  - Adopting the newly added framework monitoring memory access pattern in Linux kernel called `DAMON` to reduce the access tracking overhead

# HW-assisted tiered memory organization
(invisible DRAM to OS)

# HW-assisted tiered memory organization

- Transparent to software
  - Any software modification is not required

- We are curious about the commodity operating systems work well
  - Modern memory management is highly optimized DRAM-only systems
    - Without consideration of heterogenous (hybrid or tiered) memory systems
    - Only NUMA characteristics are considered

- We revisit the design and implementation of operating systems

# Recall memory placement: `local-first`

Let's see how the local-first policy works on the HW-assisted tiered memory system



The `local-first` placement policy leads to spending time back and forth between the local DRAM cache and the Optane main memory while the remote DRAM cache is idle

# AutoNUMA is considered harmful

Again, it is designed for DRAM-only systems



AutoNUMA balancing may degrade performance on tiered memory systems. If the local DRAM cache does not have enough space, the application can experience frequent DRAM cache misses while not utilizing the remote DRAM cache

# Our approach: `dram-first`
Exploiting such hardware characteristics in placing memory (pages)



Threads

CPU-0

Interconnection links

CPU-1

Memory Controller

Memory Controller

DRAM Cache

Allow memory allocation that can fit in DRAM cache

Make allocation on remote memory to utilize remote DRAM cache

Back to local memory

**1** Chunk #1 **3** Chunk #3

**2** Chunk #2

Main memory

Main memory

**S**

Do not always bring remote (DRAM) pages

# Preliminary results: latency & bandwidth

Experimental environments
- Intel Xeon Gold 5218 (16cores) x 2 sockets
- Two DRAM 16GB dimms and two DCPMM 128GB dimms per socket
- Linux kernel 5.3 with Ubuntu 18.04 server distribution

# Remaining challenges in tiered memory systems

1.  Demoting or migrating pages to Optane memory suffers from the <u>limited memory bandwidth</u> and leads to <u>write amplification problems</u>

    *Random write with 256B granularity

# Remaining challenges in tiered memory systems
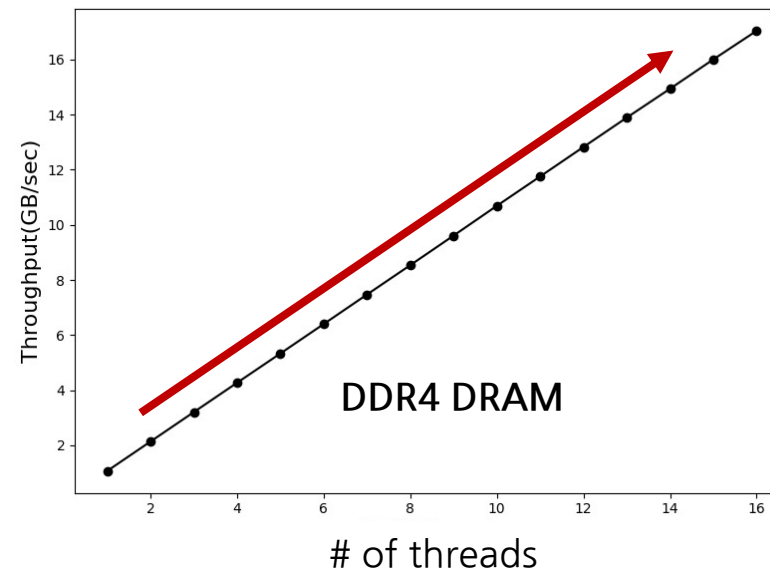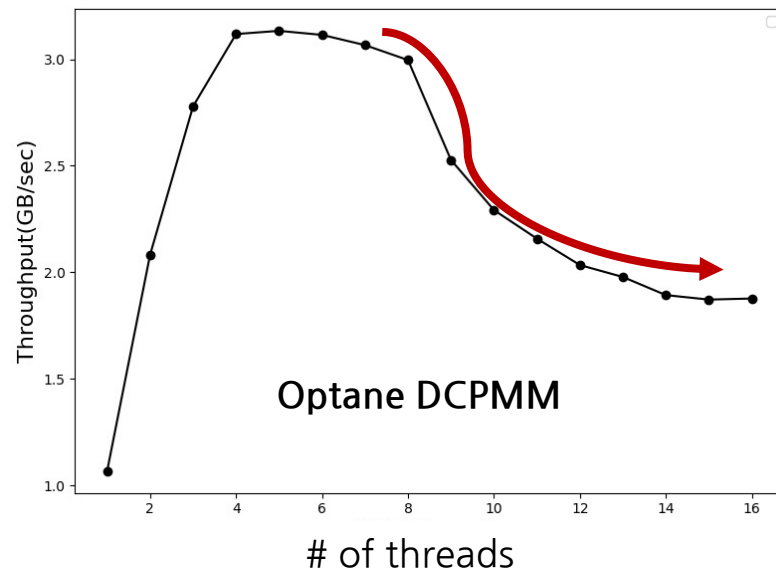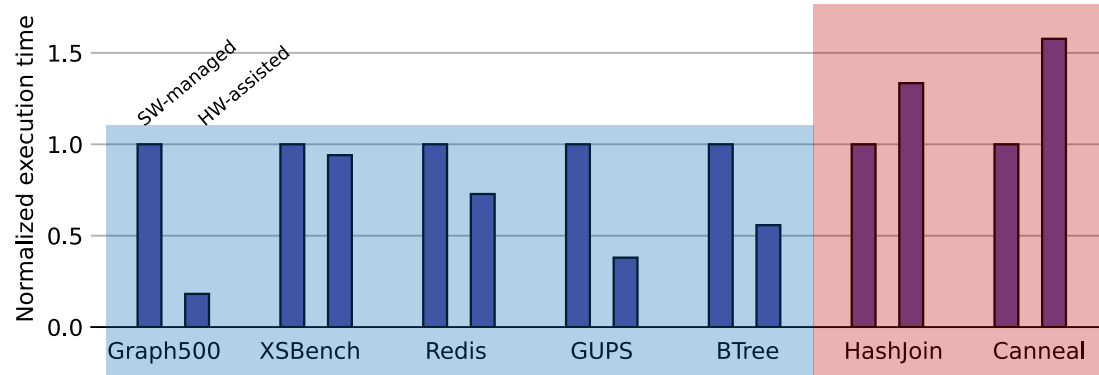
1. Demoting or migrating pages to Optane memory suffers from the <u>limited memory bandwidth</u> and leads to <u>write amplification problems</u>

2. Minimizing DRAM cache conflict misses within an Optane memory node
   - DRAM cache is organized as a <u>direct-mapped cache</u>
   - Two more memory blocks cannot be mapped to a single cache set
   - Note that the DRAM cache indexing scheme has not been disclosed

# Thank You!

Artifact available at https://github.com/csl-ajou/AutoTiering

Jeongseob Ahn

E-mail: jsahn@ajou.ac.kr
Web: https://jeongseob.github.io

AJOU UNIVERSITY