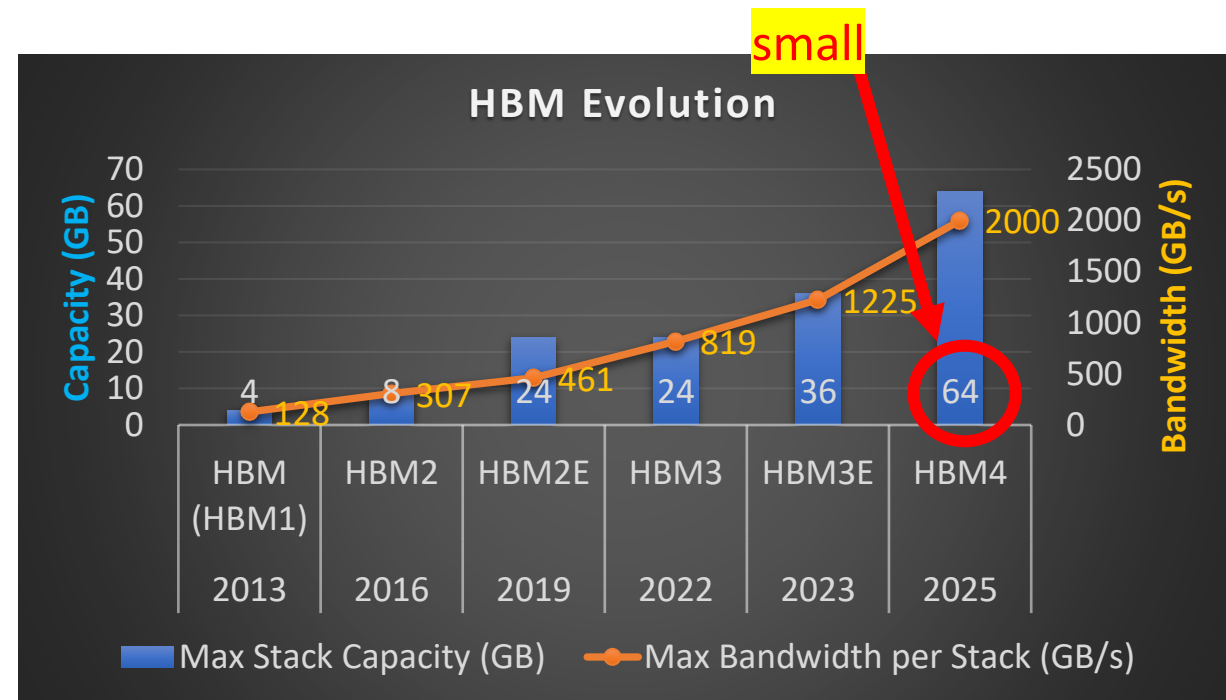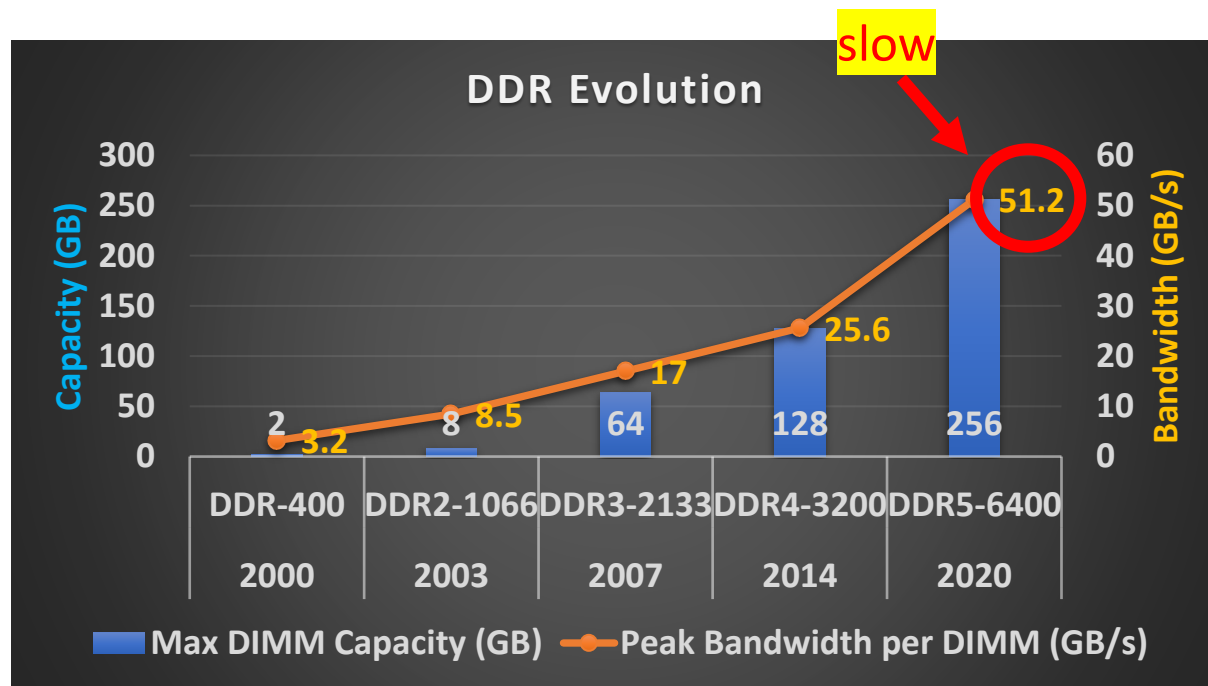# System Test Results with
# NVMe-over-CXL (NVMe-oC) Memory Mode

San Chang and Bernard Shung
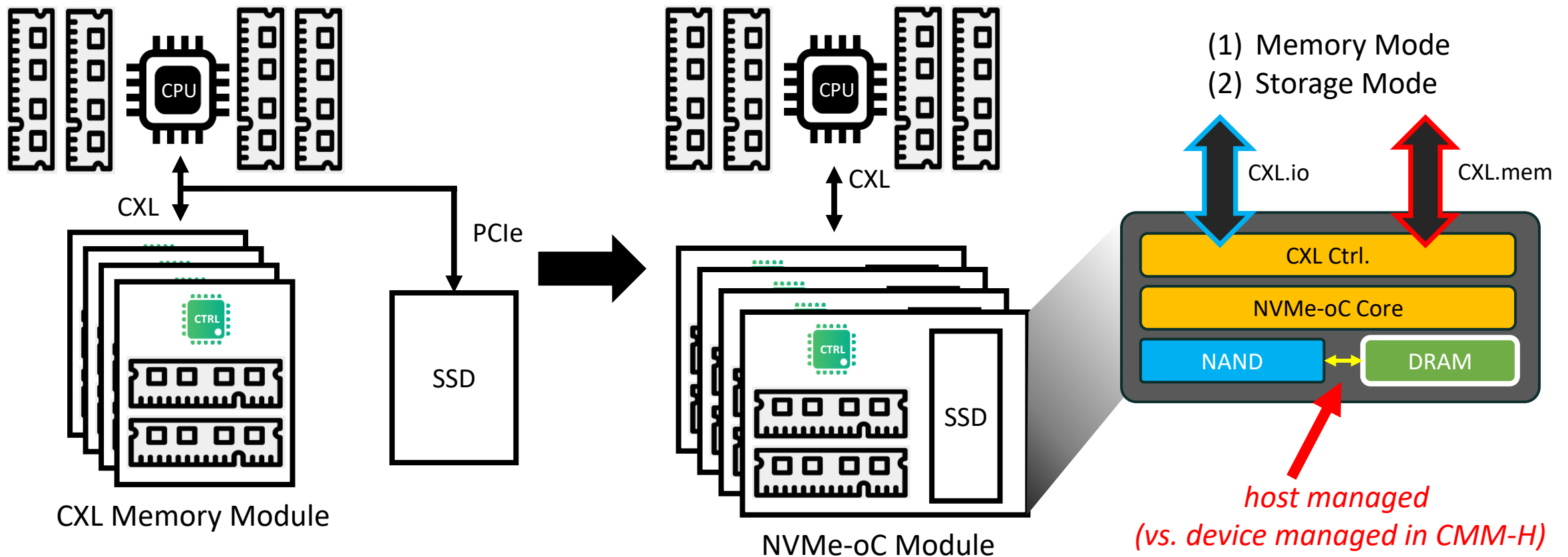Wolley Inc.

San Chang and Bernard Shung
Wolley Inc.

# Memory Is Not Scaling Fast Enough

- DRAM bandwidth grows slowly; HBM capacity gains are incremental
- Emerging workloads such as AI demand memory systems far beyond today's limits

# What is NVMe-over-CXL (NVMe-oC)?

*"Combines DRAM and NAND into a unified and CXL-attached memory module"*



CXL Memory Module

NVMe-oC Module

(1) Memory Mode
(2) Storage Mode

CXL.io

CXL.mem

CXL Ctrl.

NVMe-oC Core

NAND

DRAM

*host managed
(vs. device managed in CMM-H)*
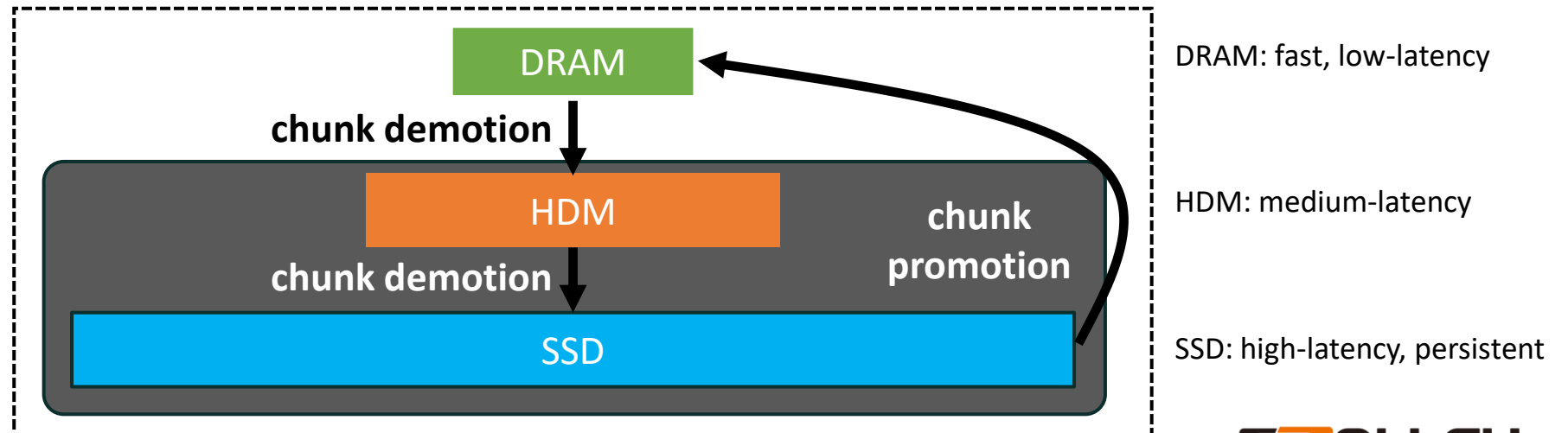
the Future of Memory and Storage

# NVMe-oC Memory Mode

- Memory Access (load/store) will either fall into either DRAM or HDM
- Hierarchical Caching with Unified Address Space
  - DRAM <-> HDM: exclusive cache pair (no duplicates, only one valid copy)
  - HDM <-> SSD: write-back strategy

**(powered by DAX-tiering)**

**Unified Memory Address Space  (up to SSD capacity)**

DRAM

**chunk demotion**

HDM

**chunk demotion**

**chunk promotion**

SSD

DRAM: fast, low-latency

HDM: medium-latency

SSD: high-latency, persistent

*the Future of Memory and Storage*

# Experiment Setup & Benchmark

| Item | Description | |
|------|-------------|---|
| CPU | Intel Xeon (Model 173), 144 cores (72 cores per socket × 2 sockets) | |
| Memory | 2 DDR5 RDIMM, 6400 MT/s, 32GB each (total bandwidth: 100GB/s) | *Capacity: 64GB* |
| NVMe-oC | PCIe Gen3 x8 (8GB/s), 16GB HDM + 64GB SSD | |
| OS | Fedora Linux 40 (Workstation Edition) | |
| kernel | 6.8.5-301.fc40.x86_64 | |

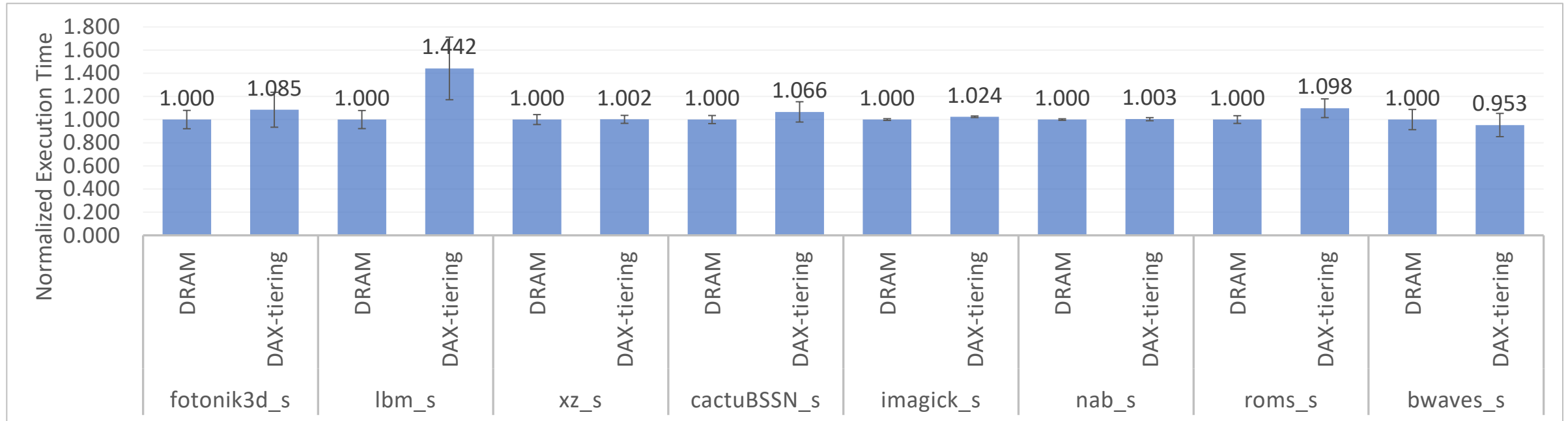### NVMe-oC Memory Mode Configurations

| | | |
|------|-------------|---|
| DAX-tiering | 16GB DRAM + 16GB HDM + 64GB SSD | *Capacity: 64GB* |

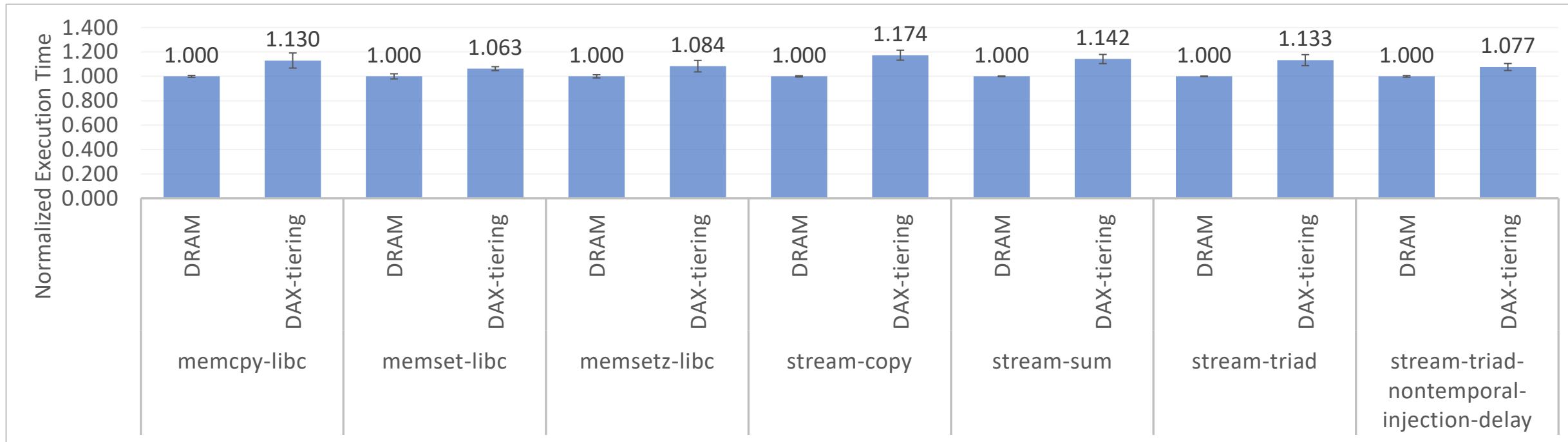| Benchmark Suite | Domain | Workload Style |
|-----------------|--------|----------------|
| SPEC 2017 | CPU & memory compute | Single/multi-core integer/FP jobs |
| multichase/multiload | Memory subsystem | STREAM-like copy/compute/write microkernel with NT stores |
| XSBench | Scientific computing / HPC | Irregular memory access patterns |
| YCSB | Cloud NoSQL / key-value stores | Configurable DB access workloads |

# SPEC 2017 – DRAM vs. NVMe-oC Memory

- NVMe-oC achieves near-DRAM performance with average CPU stall of ~1.08 (vs. DDR5: 1.0, CMM-H: ~1.7 based on Samsung paper), delivering lower stall time and outperforming CMM-H



Fits in DRAM + HDM

the **Future** of Memory and Storage

# multichase/multiload – DRAM vs. NVMe-oC Memory

Able to provide near-DRAM performance when the working set fits within DRAM + HDM



Normalized Execution Time

| | memcpy-libc | | memset-libc | | memsetz-libc | | stream-copy | | stream-sum | | stream-triad | | stream-triad-nontemporal-injection-delay | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DRAM | DAX-tiering | DRAM | DAX-tiering | DRAM | DAX-tiering | DRAM | DAX-tiering | DRAM | DAX-tiering | DRAM | DAX-tiering | DRAM | DAX-tiering |
| Value | 1.000 | 1.130 | 1.000 | 1.063 | 1.000 | 1.084 | 1.000 | 1.174 | 1.000 | 1.142 | 1.000 | 1.133 | 1.000 | 1.077 |

**Fits in DRAM + HDM**

FMS the Future of Memory and Storage

VOLLEY

# XSBench (DRAM vs. NVMe-oC DAX-tiering)

- Our NVMe-oC with DAX-tiering maintains near-DRAM performance up to workload I, while CMM-H, without host DRAM support, slows down by 1.58x to 5.56x as the memory footprint grows from 80 MB to 40 GB
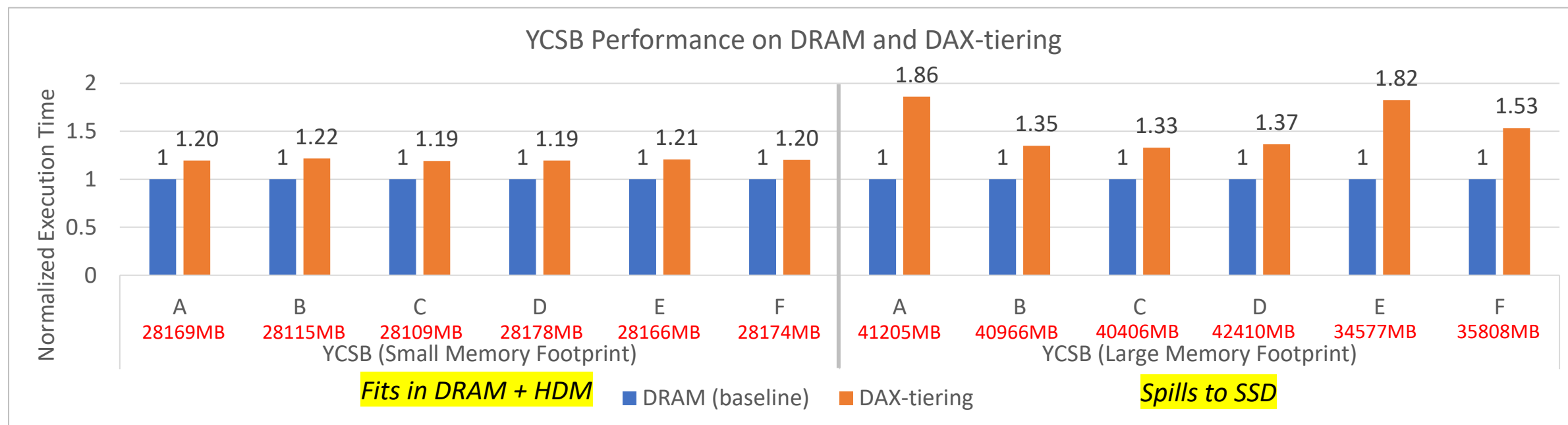


XSBench - Normalized Execution Time (DRAM vs. DAX-tiering)

Spills to SSD

Fits in DRAM + HDM

| | A (80MB) | | B (160MB) | | C (320MB) | | D (640MB) | | E (1279MB) | | F (2558MB) | | G (5117MB) | | H (10234MB) | | I (20468MB) | | J (40936MB) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DRAM | 1.000 | | 1.000 | | 1.000 | | 1.000 | | 1.000 | | 1.000 | | 1.000 | | 1.000 | | 1.000 | | 1.000 | |
| DAX-tiering | | 0.934 | | 0.937 | | 0.973 | | 0.966 | | 0.985 | | 0.985 | | 1.209 | | 1.000 | | 1.203 | | 6.531 |

FMS the Future of Memory and Storage

VOLLEY

# YCSB Performance
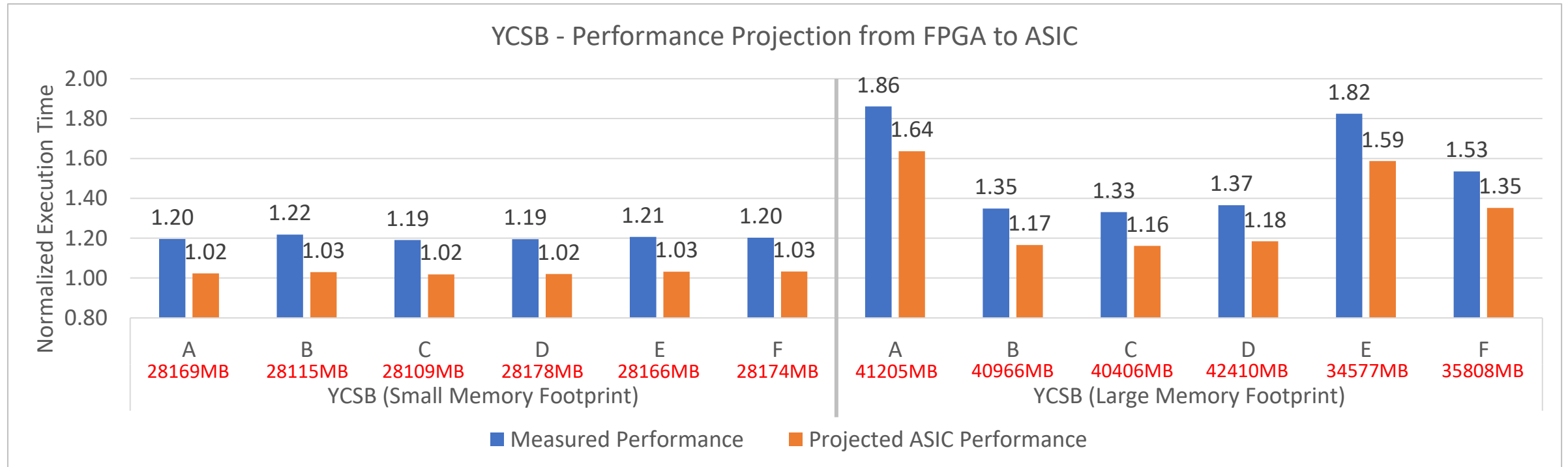## (DDR5 vs. NVMe-oC DAX-Tiering )

- **DAX-tiering delivers 83% and 65% of DDR5 performance** under small and large memory footprints respectively, while costing less than half of DDR5. This results in **performance-per-dollar gains of 67% (small footprint) and 30% (large footprint)**
  - Small memory footprint workloads (fits in DRAM + HDM) show an average of 1.2× slowdown
  - Large memory footprint workloads (spill over to SSD) show an average of 1.54× slowdown



YCSB Performance on DRAM and DAX-tiering

# YCSB Performance Projection from FPGA to ASIC

- With ASIC, **DAX-tiering delivers 97% (vs. 83%) and 74% (vs. 65%) of DDR5 performance** under small and large footprints, at less than half the cost. This leads to the **performance-per-dollar gains of 94% (small footprint) and 48% (large footprint)**



YCSB - Performance Projection from FPGA to ASIC

# Performance Comparison (NVMe-oC vs. CMM-H)

- NVMe-oC (with 16GB DDR4-2000)
  - (SPEC CPU2017) NVMe-oC, which combines host DRAM and device DRAM as cache, shows an average 8.4% slowdown compared to DDR5-L
  - (XSBench) Performance degrades by 1.0× to 6.5× as the memory footprint increases from 80 MB to 40 GB

- CMM-H (with 16GB DDR4-2666)
  - (SPEC CPU2017) shows an average 70% performance degradation compared to DDR5-L (local DRAM)
  - (XSBench) Performance degrades by 1.58× to 5.56× as the memory footprint increases from 80 MB to 40 GB

Zeng, Jianping, et al. "Performance Characterizations and Usage Guidelines of Samsung CXL Memory Module Hybrid Prototype." *arXiv preprint arXiv:2503.22017* (2025).

# Feature Comparison (NVMe-oC vs. CMM-H)

| | CMM-H | NVMe-oC |
|---|---|---|
| Host Interface | PCIe Gen5x8 | PCIe Gen5x8 |
| Media | DRAM + SSD | DRAM + SSD |
| Memory Mode Capacity | Up to SSD total capacity | Up to SSD total capacity |
| Memory Mode Namespace | Single | Multiple (user configurable) ✅ |
| HDM Management in Memory Mode | Device Cache Controller | Host CPU |
| Dual Mode (Memory/Storage) Online | NA | Support ✅ |
| Insert Application Intelligence | Limited | Support ✅ |
| Direct SSD Data to Host DRAM (Hot Data) | NA | Support ✅ |
| Dual-cache (Host DRAM + Device DRAM) | NA | Support ✅ |
| Data Prefetching | Device Cache Controller | Support ✅ |
| Hot/Cold Data Detection | Device Cache Controller | Driver (TRACE) |
| Cache Policy | 8-way, 4KB management, LRU, MRU (HW implementation) | N-way, M KB management, LRU, MRU, etc., (software defined) |

the **Future** of **Memory** and **Storage**

# Takeaways

- NVMe-oC memory mode enables cost-effective memory expansion over CXL for data-intensive workloads
  - Jointly manages host DRAM, device DRAM (HDM) and NAND with flexible caching
  - No code changes to the existing applications
- NVMe-oC memory mode value proposition
  - Provides near-DRAM performance when the working set fits within DRAM plus HDM, serving as a standard CXL memory module for memory expansion
  - Delivers better performance/$ (30~90% improvement) when the working set exceeds DRAM plus HDM, outperforming CMM-H through host-managed caching intelligence

San Chang, san@wolleytech.com

the Future of Memory and Storage