# How CXL Will Change
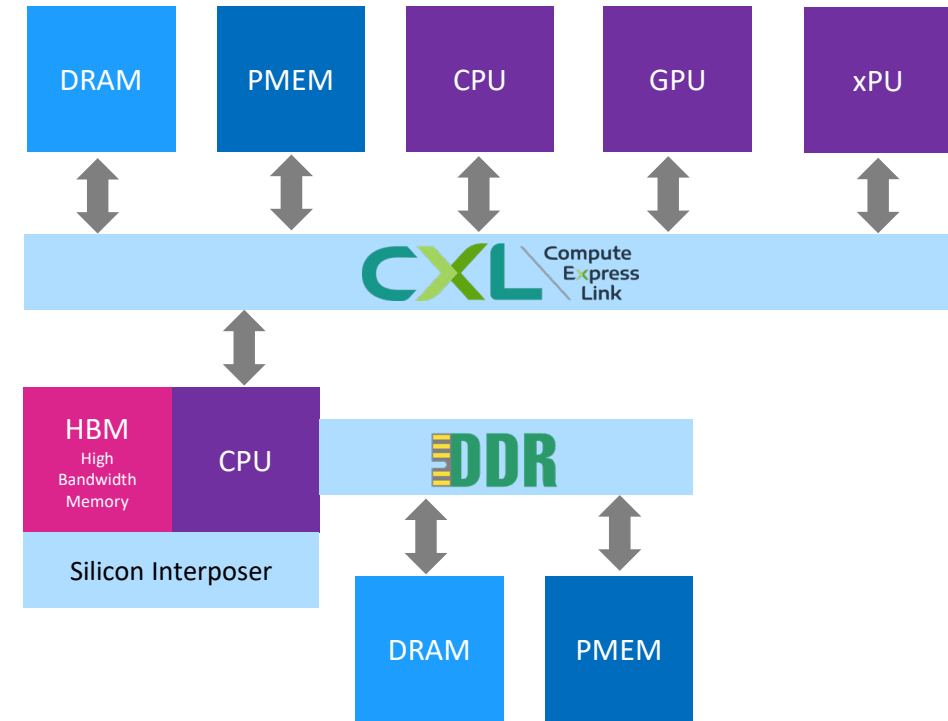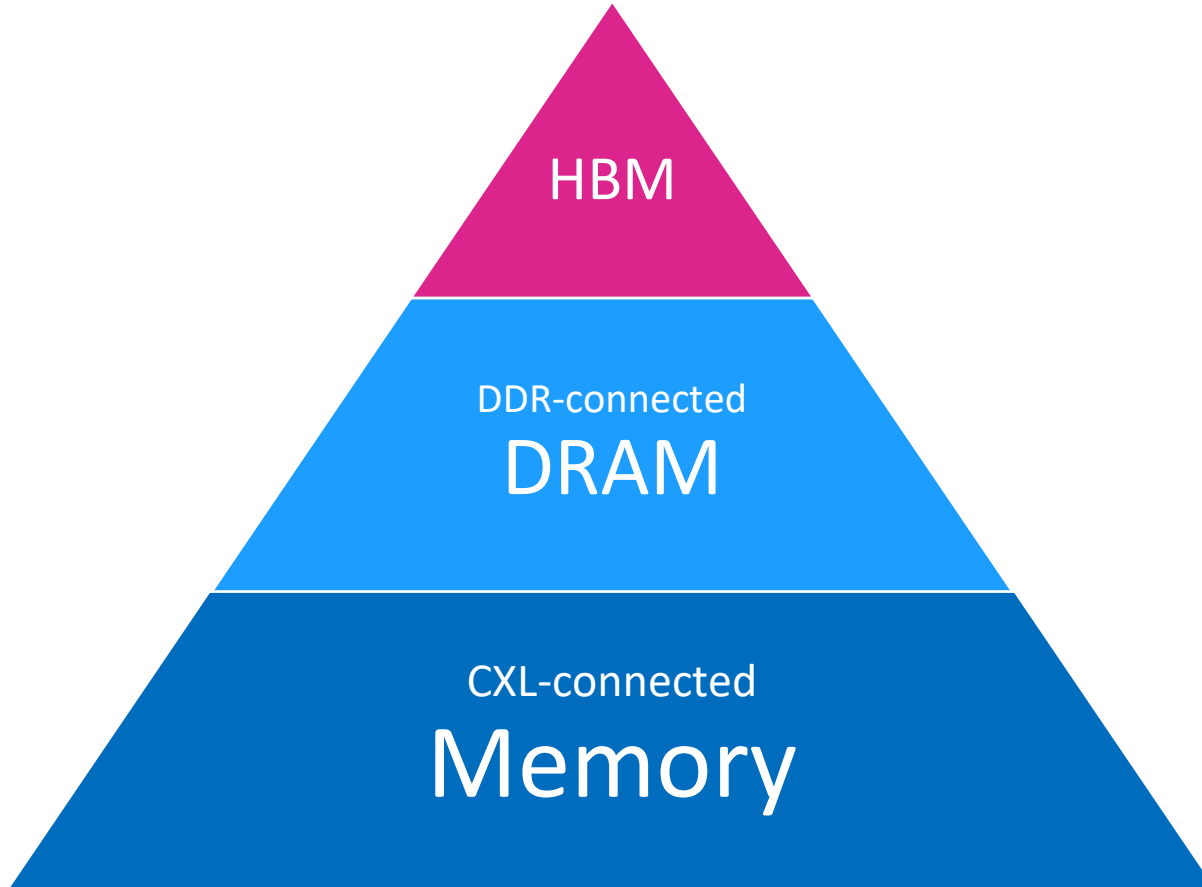# The Datacenter
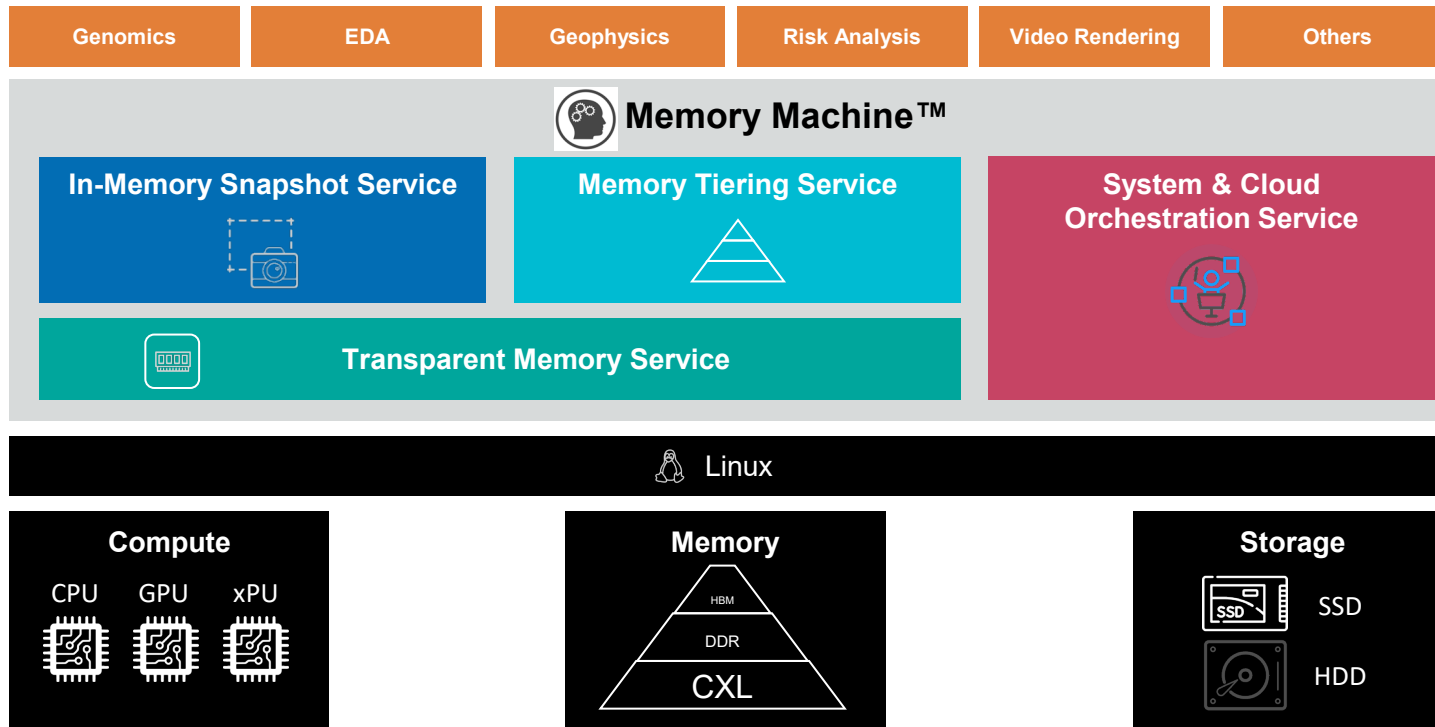
Presented by Bernie Wu, VP Business Development

bernie.wu@memverge.com

# CXL – A game changer



HBM

DDR-connected
DRAM

CXL-connected
Memory

DRAM | PMEM | CPU | GPU | xPU

CXL Compute Express Link

HBM High Bandwidth Memory | CPU | DDR

Silicon Interposer

DRAM | PMEM

But hardware Innovation Alone
Is Not Enough!

SNIA PERSISTENT MEMORY + SUMMIT 2022 COMPUTATIONAL STORAGE

# MemVerge Memory Machine™

| Genomics | EDA | Geophysics | Risk Analysis | Video Rendering | Others |
|----------|-----|------------|---------------|-----------------|--------|

### 🧠 Memory Machine™

| **In-Memory Snapshot Service** | **Memory Tiering Service** | **System & Cloud Orchestration Service** |
|---|---|---|

**Transparent Memory Service**

🐧 Linux

**Compute**
CPU  GPU  xPU

**Memory**
HBM
DDR
CXL

**Storage**
SSD
HDD

---

**Memory Virtualization**

- Software-defined Memory Pool

- Intelligent Auto-tiering

- **Delivers big memory capacity without application change**

**Memory Snapshot Service**

- Fully captures running application state

- An Application can be rolled back, restored or cloned from anywhere at any time

- **Delivers higher application mobility and availability**

**Memory Object Sharing**

- Enables Objects to be stored/shared in-memory with tiering into storage

# Use cases



- Genomics — secondary and tertiary analytics sequencing/assembly
- EDA/CAE — simulation and modeling
- AI/ML/Analytics — Deep learning training, batch and interactive IDE
- Media & Entertainment — Simulation/Rendering & SFX
- HPC- simulation/modeling
- FSI- analytics/decision support, low-latency persistent messaging
- In-Memory Databases — node consolidation, performance improvements, backup/recovery
- Cloud- KVM snapshots, Spot-instance for non-FT workloads

PERSISTENT MEMORY
+ SUMMIT 2022
SNIA COMPUTATIONAL STORAGE
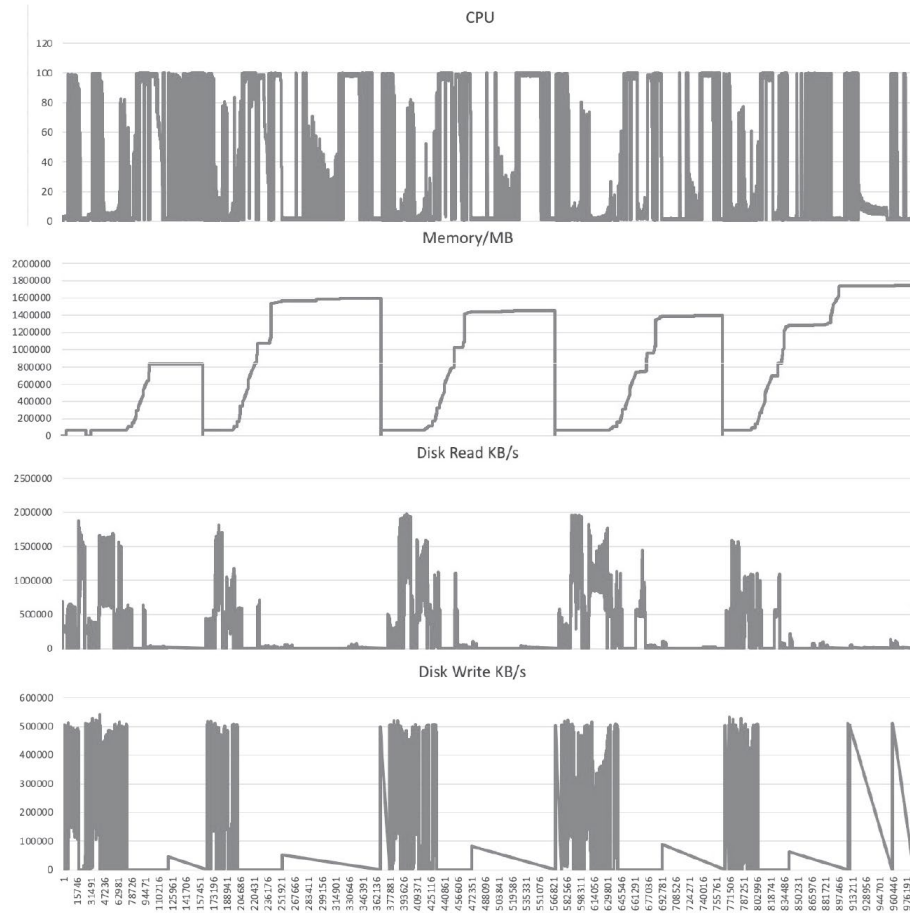
# CXL Benefits

- Industry-wide Standard that will initially benefit Memory-Bound applications
    - Increase total amount and types of memory
    - Memory Bandwidth
    - HPC and AI/ML type workloads biggest benefactors

- Cloud
    - Ability to Improve overall utilization, efficiency of heterogeneous computing
- Composability (HW-level)
    - Memory becomes a First-Class Citizen, eventually consolidated into centralized pools
- Collaboration
    - Facilitate Data-sharing across applications and servers
- Concurrency
    - Enable concurrent heterogenous processes to run on same in-memory data sets to reduce wall-clock time
- Cache
    - Reduced movement of data to/from storage and memory as a result of larger pools/caches including persistence

PERSISTENT MEMORY + SUMMIT 2022 COMPUTATIONAL STORAGE

# Challenges to CXL Adoption

- Many HPC and AI/ML application/pipelines are "pets" with wide-ranging memory access patterns and behaviors.

# HPC "Pet" Workload Compute Profile Example



MetaSpades Application alternates between compute, memory, and IO bound

https://www.biorxiv.org/content/10.1101/2022.04.20.488965v1.full.pdf

PERSISTENT MEMORY
+
SNIA SUMMIT 2022
COMPUTATIONAL STORAGE

# First Step in Adopting CXL: Achieving Application Transparency

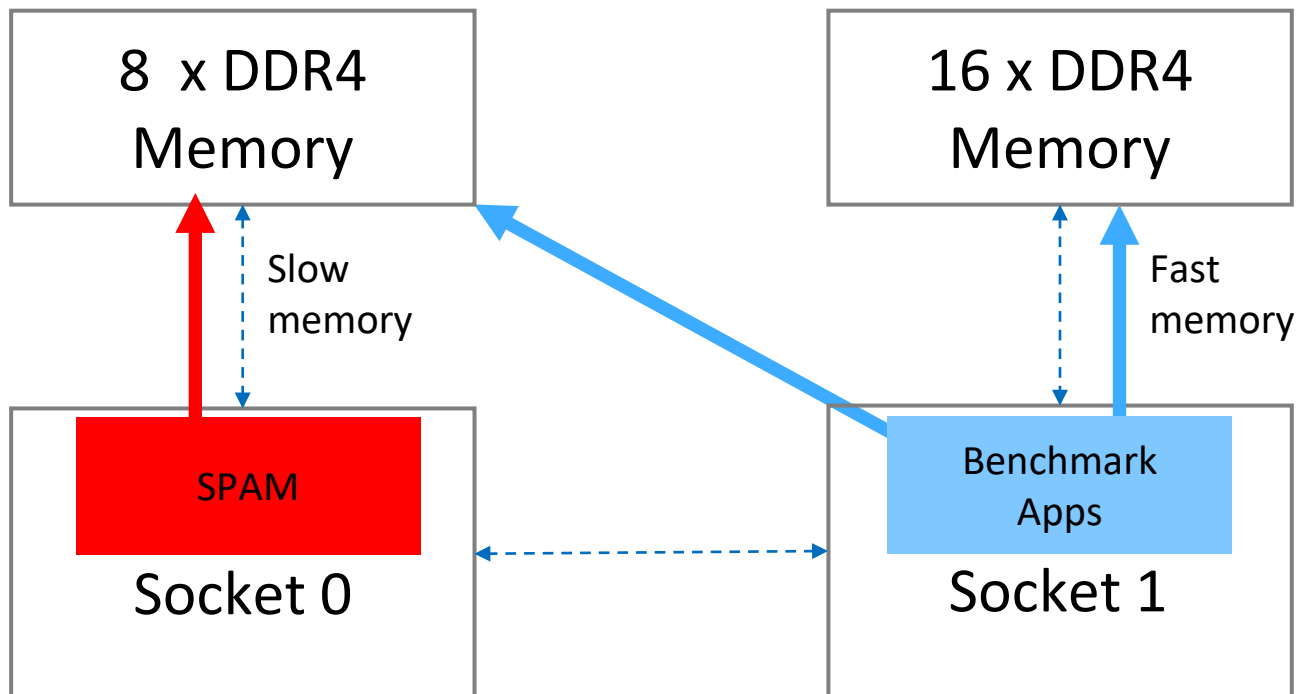Goal: Targeted HPC and AI/ML Applications Must run Better, Faster, Cheaper without Modification

Pre-CXL objective: Transparent, Automated Memory Tiering

      A). Memory Capacity: Memory tiering between DRAM and PMEM on same DDR bus

      B) Memory Bandwidth/Latency: Tiering across local and remote NUMA DRAM. CXL 1.1 memory expansion latency expected to equivalent to one NUMA hop

Validate impact of various pooling, profiling, policy, and placement approaches to memory tiering

# Emulating CXL Two-Tiered Memory with DRAM

- DRAM 384GB DDR4 in total
  - Slow tier on socket 0:  DDR4 2666 MHz 16GB x 8
    - Use less DIMMs and running spam app to further consume memory bandwidth
  - Fast tier on socket 1:  DDR4 2666 MHz 16GB x 16
    - Benchmarking app runs on socket 1, has the fastest access to its local memory
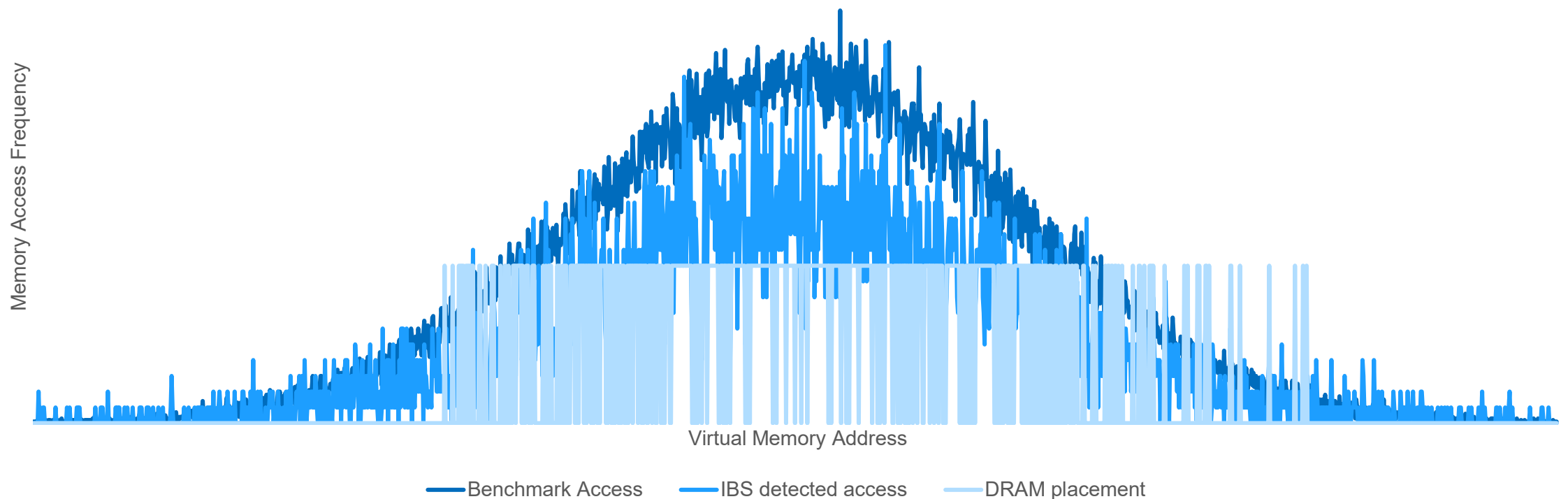


|  | Fast Tier | Slow Tier |
| --- | --- | --- |
| Access Latency | 108ns | 335ns |

*Measured with memory latency checker

9

MEMVERGE CONFIDENTIAL

PERSISTENT MEMORY
+ SUMMIT 2022
COMPUTATIONAL STORAGE

# Memory Profiling Supplemented by HW-based IBS

- MemVerge synthetic memory load generator simulates memory access that follows a Gaussian distribution (i.e., data locality)
- Memory Machine places faster DRAM to back the hottest virtual memory addresses detected by IBS-assisted memory profiler



Memory Access Frequency / Virtual Memory Address

—— Benchmark Access  —— IBS detected access  —— DRAM placement

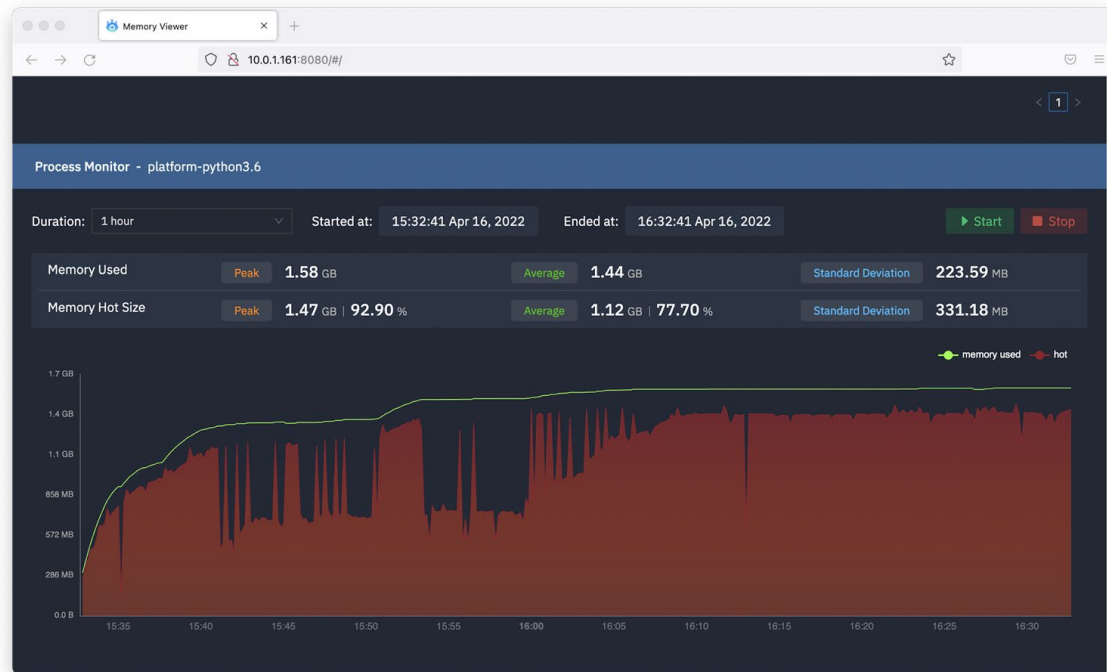SNIA + SUMMIT 2022 COMPUTATIONAL STORAGE

# Synthetic Benchmark Performance

- MemVerge synthetic memory load generator simulates memory access that follows a Gaussian distribution (i.e., data locality)
  - Benchmark used 2GB fast memory + 6GB slow memory
  - Measured memory access bandwidth

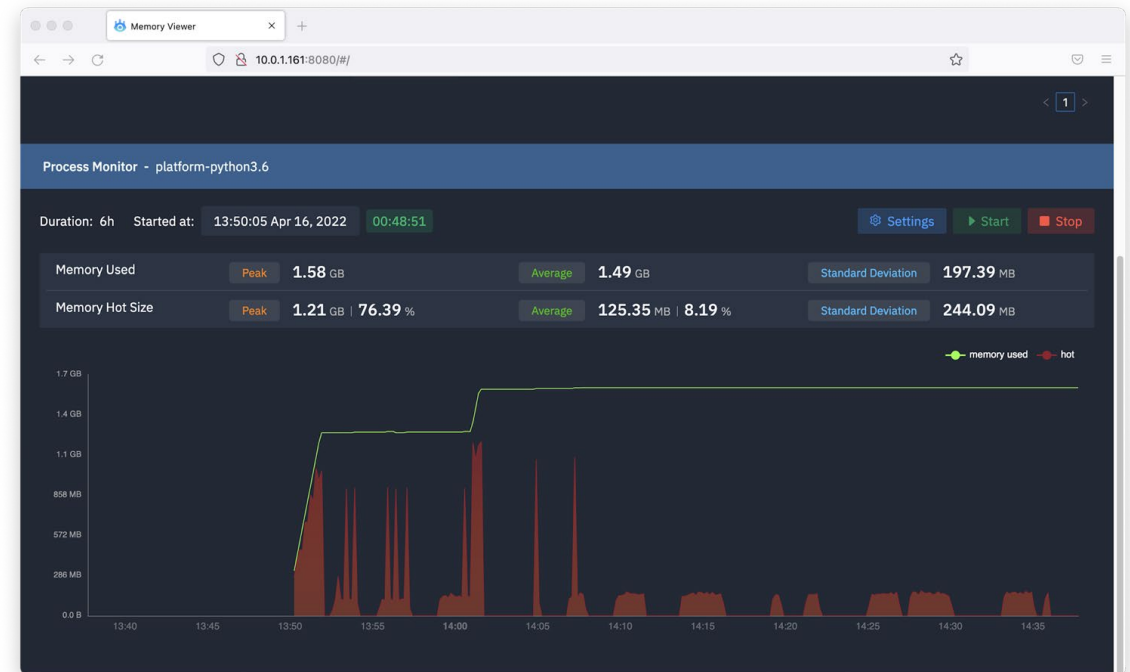Bandwidth (GB/Sec)



| | |
|---|---|
| 7 | 7.7 |
| Memory Machine with Default Profiler | Memory Machine with AMD IBS Profiler |

# Memory Viewer for Looking at Application-Level Memory Heatmap

**Not suited for Memory Tiering**

**Ideal for Tiering**

# Deployment Scenarios for Bare Metal, Kubernetes, and KVM

| | | | |
|---|---|---|---|
| Application (×3) | Application (×4) | Application (×4) | Application (×3) |
| Memory Machine (as DaemonSet) | Guest OS (Linux) / Guest OS (Win) | Memory Machine / Memory Machine | MemoryMachine (as DaemonSet) |
| Kubernetes Container | KVM / KVM | Guest OS (Linux) / Guest OS (Linux) | Kubernetes Container |
| Linux Host OS | Memory Machine | KVM / KVM | Guest OS (Linux) |
| Bare-metal Server w/CXL | Linux Host OS | Linux Host OS | KVM / KVM |
| | Bare-metal Server w/CXL | Bare-metal Server w/CXL | Linux Host OS |
| | | | Bare-metal Server w/CXL |

Flexibility of Memory Machine deployment allows operators to decide level of granularity of memory services and/or target specific applications for optimization

SNIA PERSISTENT MEMORY + SUMMIT 2022 COMPUTATIONAL STORAGE

# Summary

- CXL will help revolutionize datacenter architecture starting with memory
  - Memory as first-class citizen – more bandwidth, varying latencies, varieties, and degrees of composability

- CXL will help accelerate many HPC & AI-ML application workflows – faster time to insight

- Software-defined Memory is key to provisioning and managing various CXL memory pools that optimize application/workflow pipeline behavior.

- Partner with us in Shaping the Future of Big Memory!

# Thank You