# I/O Coordination for Better Resource Sharing

## From HPC to AI Storage
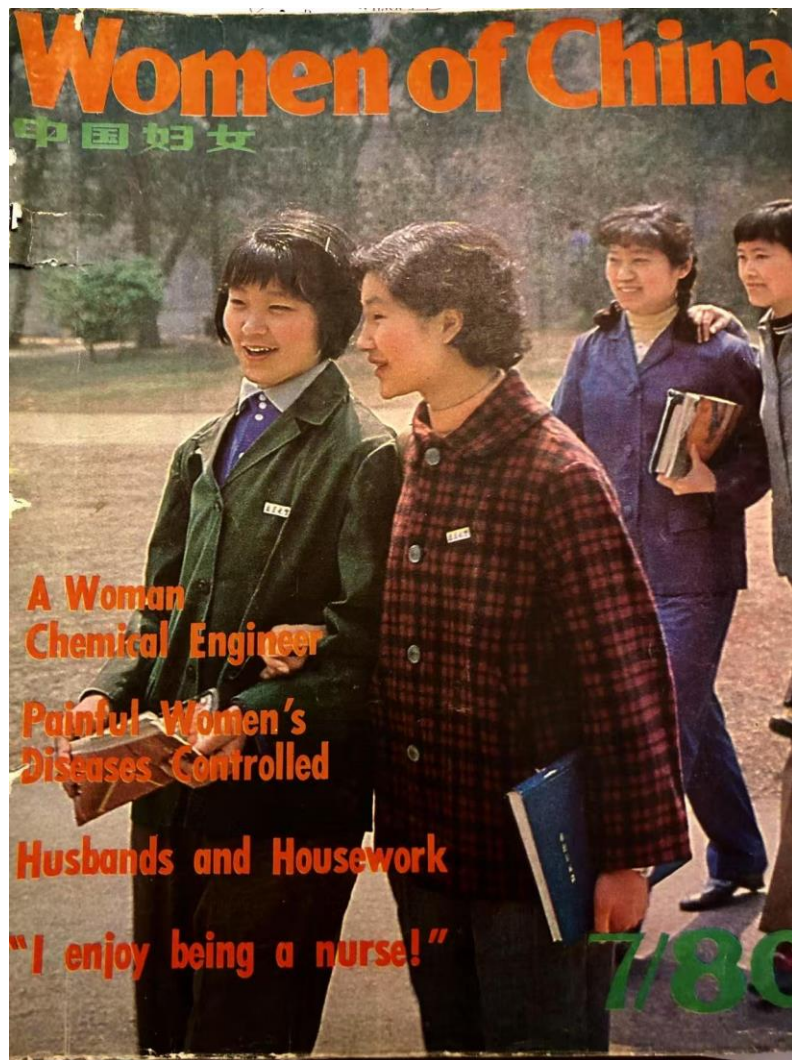
**Xiaosong Ma**

Department of Computer Science, MBZUAI

*HotStorage 2025*

# Compute Resources Shared Then



**First data center – 1950s**

Source: https://opticalcloudinfra.com/index.php/what-why-and-how/short-data-center-history
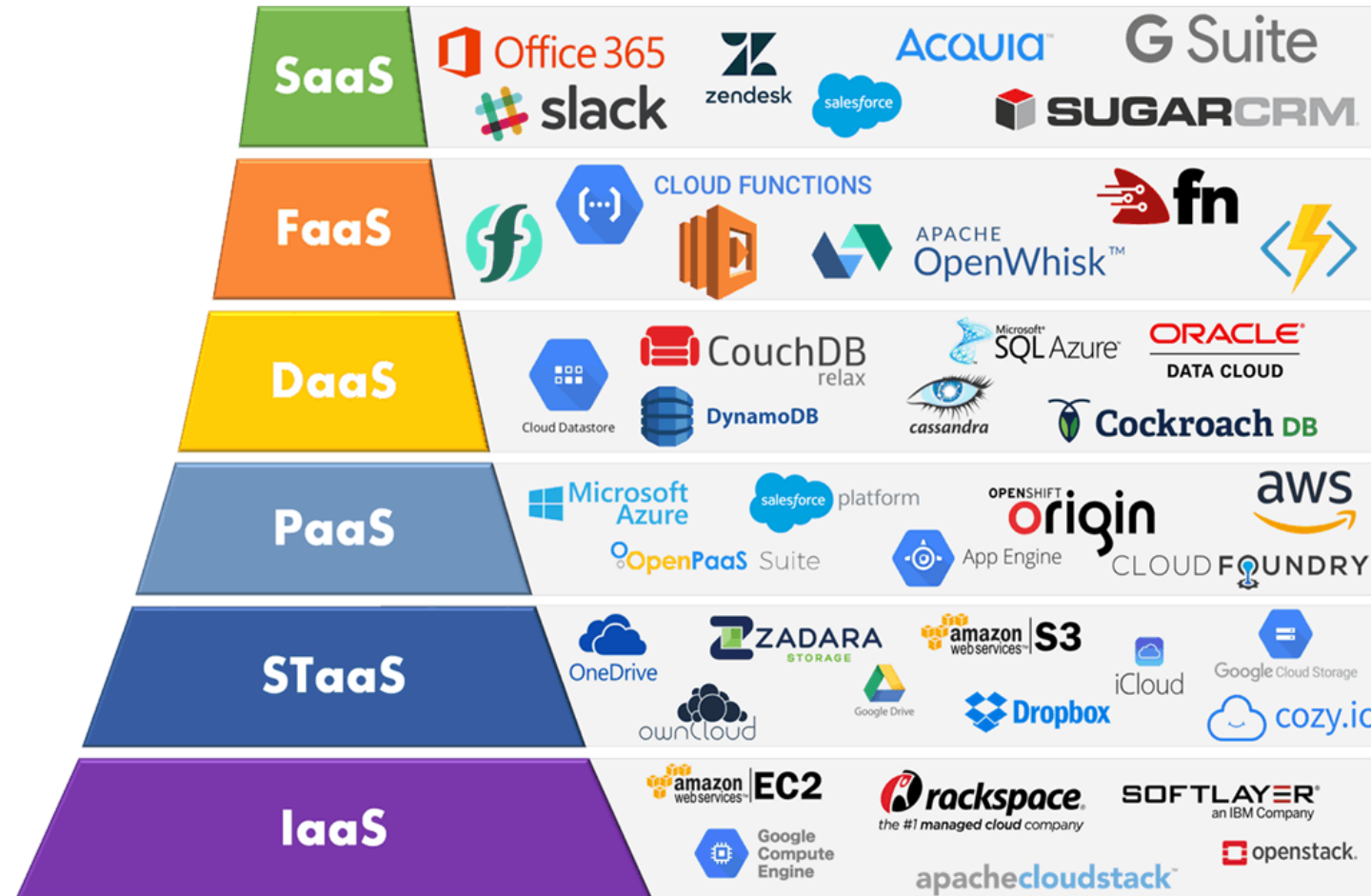
# And Again Now



**Supercomputers**
Source: Oak Ridge National Lab



**Data centers**
Source: https://www.crn.com/news/data-center/google-unveils-new-750m-data-center-as-part-of-9-5b-goal



**Cloud platforms and services**
Source: https://kinsta.com/blog/go`ogle-cloud-vs-aws/

4

# **Sharing-Friendly Hardware Platforms**

❑ Multi-core processors

  ➢ Dozens of cores

  ➢ CPU caches optimized for multi-tenancy (e.g., large L2 caches)

  ➢ TBs of DRAM space

  ➢ Mechanisms for inter-core resource allocation (cache ways, memory BW)

❑ Powerful interconnect

  ➢ Fast network connections (e.g., up to 100Gbps at AWS)

  ➢ Smart NICs/DPUs offloading computation tasks

❑ High-capacity storage

  ➢ NVMe SSDs offering space, bandwidth, and IOPS for sharing

**AMD EPYC 9654**

**Source: AMD**

# Implications of Pervasive Resource Sharing

❑ Programs to run
  ➢ on unknown/changing hardware
  ➢ with unknown/changing neighbors
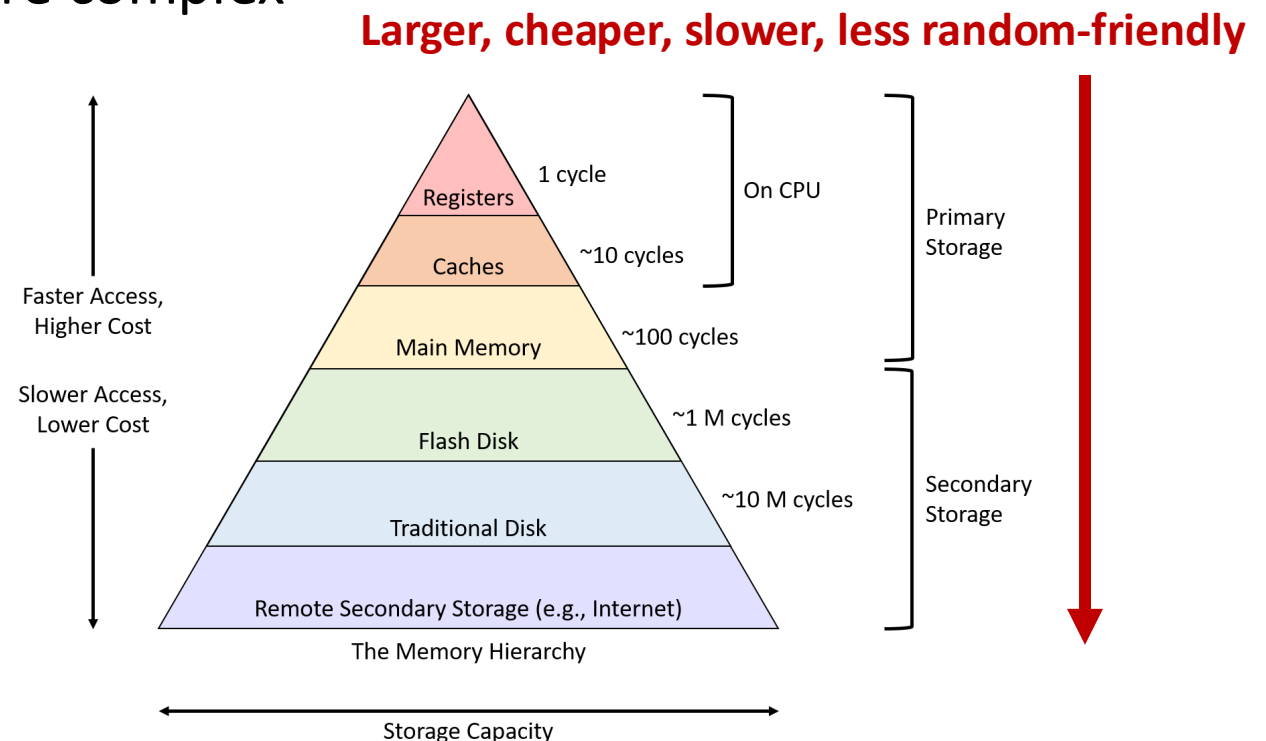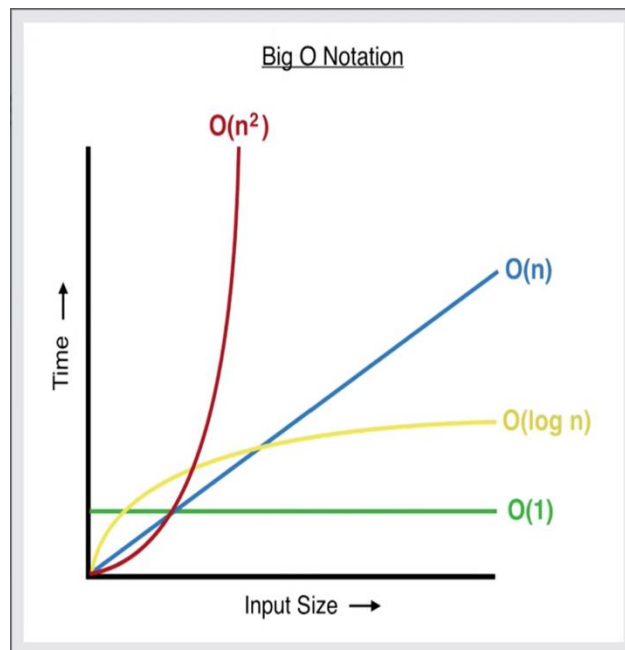
❑ Major challenge: performance portability

❑ Desirable for applications/services to
  ➢ **retain (optimized) performance across platforms**
  ➢ **achieve hardware potential**

# Challenge Lies in Storage Hierarchy

❑ Computation logic more "portable"

  ➢ Instruction execution easier in scheduling and isolation

  ➢ Current server processors w. sizable per-core resources (e.g., L1+L2 cache)

❑ Data access path deeper and more complex

**Larger, cheaper, slower, less random-friendly**



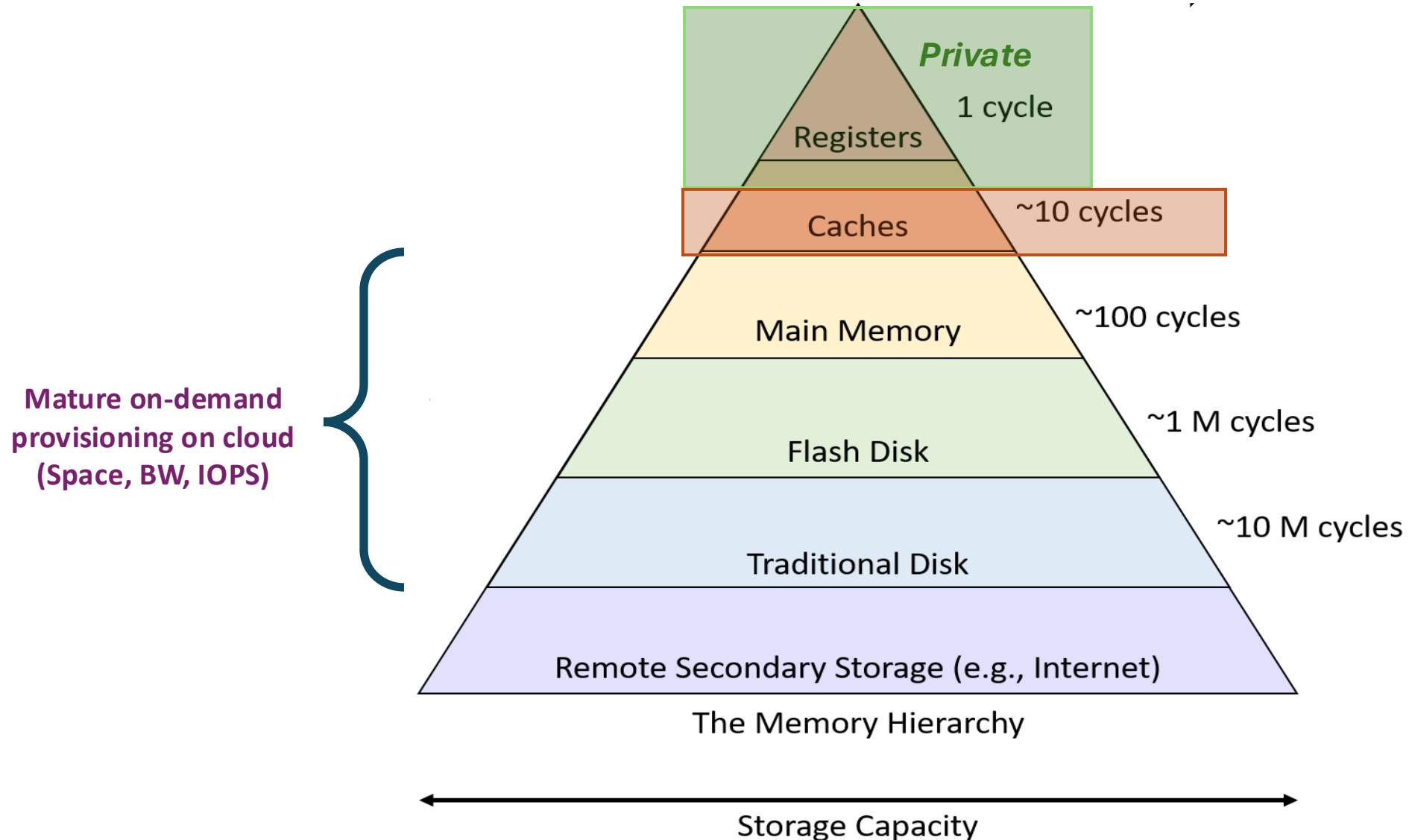Big O Notation



The Memory Hierarchy

Storage Capacity

# Storage and I/O Not Efficiently Shared

❑ Major factor leading to storage frustrations and wastes

➢ From single server to data-centers/supercomputers

❑ In this talk

➢ Sample related problems and solutions in our past research

➢ Extended I/O hierarchy for AI workloads

# Challenging Layer 1: Shared Cache



**Private**

1 cycle

Registers

~10 cycles

Caches

~100 cycles

Main Memory

~1 M cycles

Flash Disk

~10 M cycles

Traditional Disk

Remote Secondary Storage (e.g., Internet)

The Memory Hierarchy

**Mature on-demand provisioning on cloud (Space, BW, IOPS)**

Storage Capacity
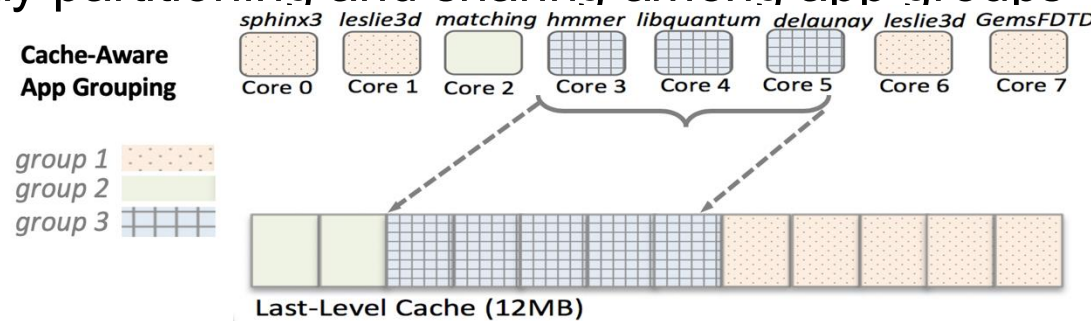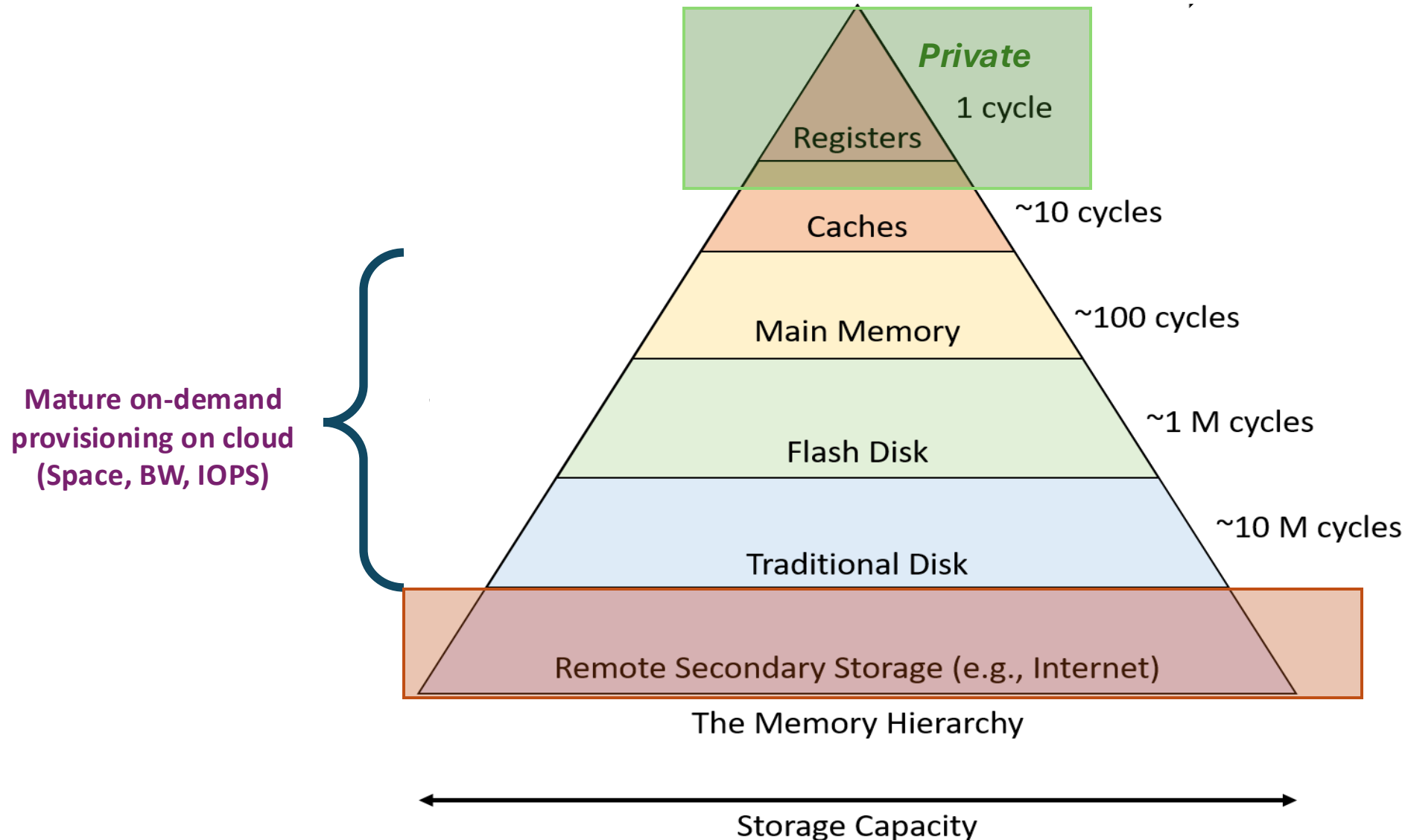
# Cache Partitioning in Commodity Multicores

- ✓ Current processors offer hardware cache-partitioning support (Intel CAT)

- Partitioning last-level cache among co-running apps
  reduces interference ➔ improves system performance

- *Kpart* [HPCA '18]

- ✗ Two key challenges limit usability of CAT
  - Current hardware implements coarse-grained way-partitioning
    ➔ hurts system performance!
  - Lacks hardware monitoring units to collect cache-profiling data

- Solution: hybrid way partitioning and sharing among app groups



- Significant performance gain on real hardware (avg 24%, max 79%)

*Private*

1 cycle

Registers

Caches — ~10 cycles

Main Memory — ~100 cycles

Flash Disk — ~1 M cycles

Traditional Disk — ~10 M cycles

Remote Secondary Storage (e.g., Internet)

The Memory Hierarchy

**Mature on-demand provisioning on cloud (Space, BW, IOPS)**

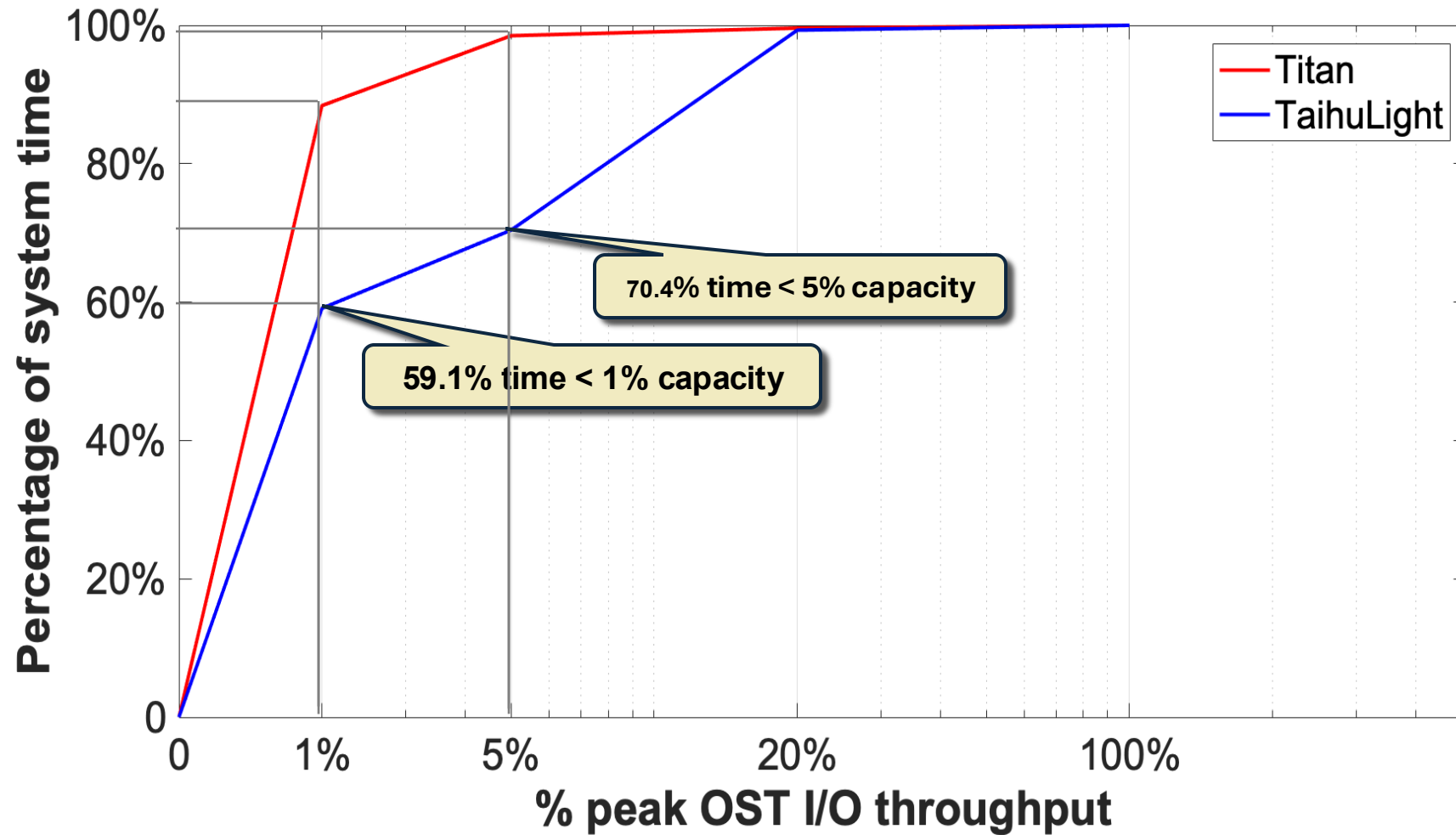Storage Capacity

# Supercomputer Storage: Not Better Shared

❑ Observation from supercomputer I/O profiling ([FAST14], [SC16], [NSDI19])

  ➢ Significant inter-job interference -> inconsistent I/O performance

  ➢ Vast majority of HPC jobs non-I/O intensive -> overall low I/O resource utilization

| Name | Value |
|------|-------|
| Total number of logged jobs | 181,969 |
| Unique applications identified | 9,998 |
| Initial I/O-intensive candidates | 95 |
| Candidates passing scope checking | 67 |
| Candidates passing minimum support | **42** |
| User-verified candidates | 8 |

**Job I/O statistics based on ORNL Titan supercomputer, 2015**

# Bottleneck, Contention Point, and Under-Utilization

# End-to-end Supercomputer I/O Monitoring [NSDI19]

❑ Understand HPC I/O for designing future systems/applications
  ➢ Lightweight end-to-end I/O resource monitoring

❑ Deployed at TaihuLight
  ➢ No user effort required
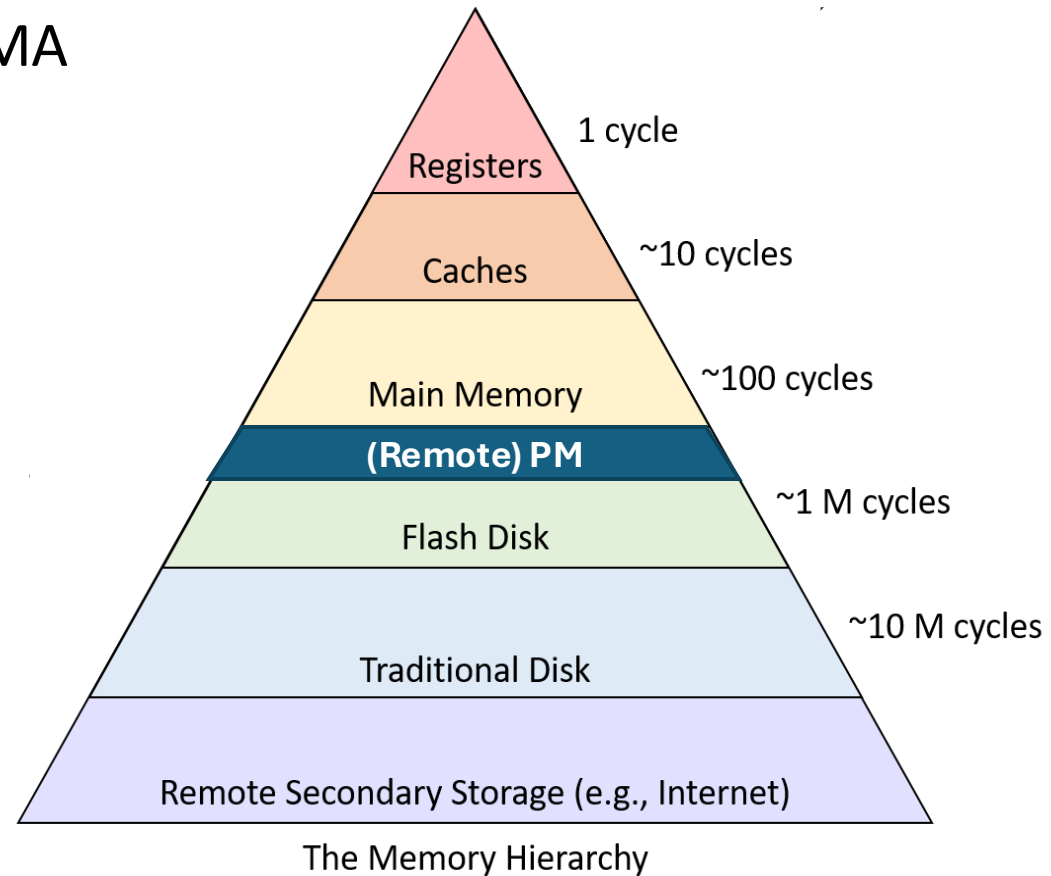  ➢ Code and monitoring data released: ***https://github.com/Beaconsys/Beacon***

❑ Findings based on 18-month monitoring on production platform
  ➢ Wide-spread adoption of inefficient I/O modes
    • Lowing both **application performance** and **hardware utilization**
  ➢ System anomalies and their behaviors (echoing findings from datacenters)
  ➢ Obscure design/configuration problems, e.g., forwarding layer cache thrashing
  ➢ Significant forwarding node load imbalance => Application-aware I/O forwarding [FAST19]

# Another Layer in Pyramid: Remote PM

❑ Persistent memory disaggregation

➢ Faster than local SSDs with RDMA

➢ Enables

- large memory buffer
- lean compute nodes

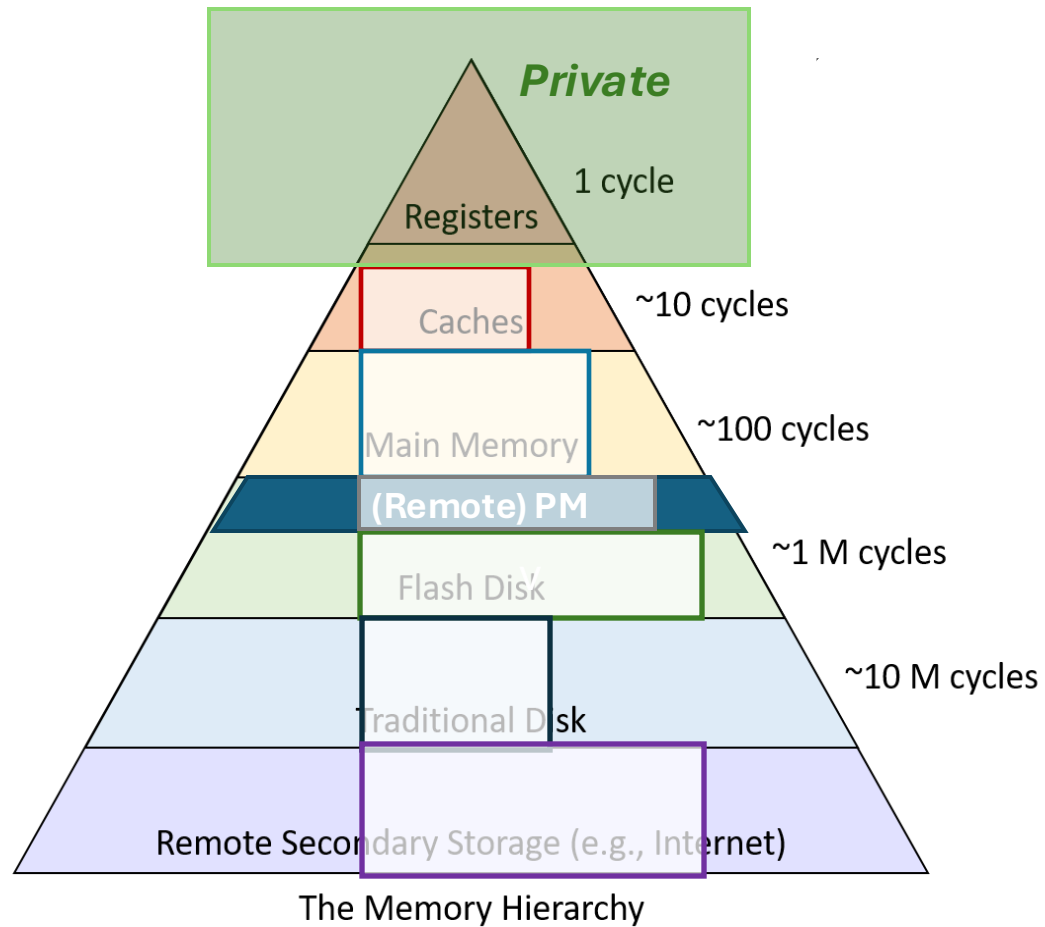| Layer | Cycles |
|-------|--------|
| Registers | 1 cycle |
| Caches | ~10 cycles |
| Main Memory | ~100 cycles |
| **(Remote) PM** | |
| Flash Disk | ~1 M cycles |
| Traditional Disk | ~10 M cycles |
| Remote Secondary Storage (e.g., Internet) | |

The Memory Hierarchy

# Cloud-Native DB on Disaggregated PM

- Distributed RDBMS [ASPLOS23, VLDB25]

**SQL Execution Engine**

**CPU Layer**

**Compute Nodes (CNs)**

SQL

TXN

**DRAM**

**(10GB)**

**Thin-CN pros**
- Core utilization
- Fast recovery

**Fast Interconnect**

**Buffer Pool**

**PM Nodes (PMNs)**

**PM (buffer and logs)**

**(2TB)**

**DRAM limitations**
- Expensive
- Low density
- Volatile

**HDD/ SDD Layer**

**Storage**

**KIOXIA XL-FLASH**

**Optane PM**          **XL-FLASH**

**All-around features:**
- **Performance close DRAM**
- **Yet cheaper**
- **High density**
- **Byte-addressable**
- **Non-volatile**

# "Storage Container" All the Way



The Memory Hierarchy

- Private
- 1 cycle — Registers
- ~10 cycles — Caches
- ~100 cycles — Main Memory
- (Remote) PM
- ~1 M cycles — Flash Disk
- ~10 M cycles — Traditional Disk
- Remote Secondary Storage (e.g., Internet)

**Within single workload: how to better use allocated hardware?**

# Sample Memory/Cache Performance

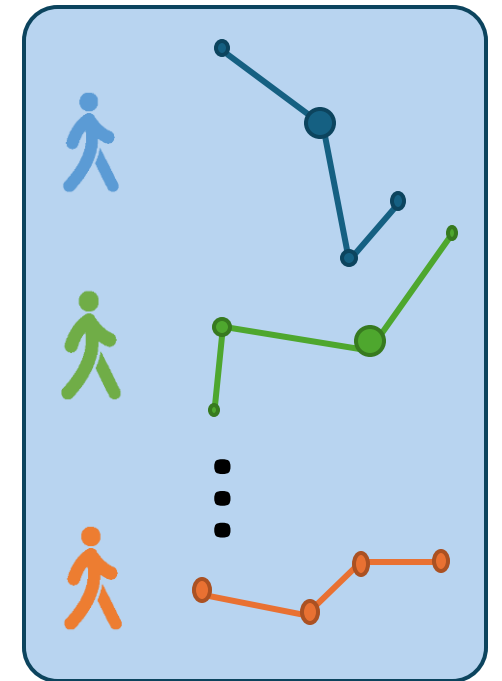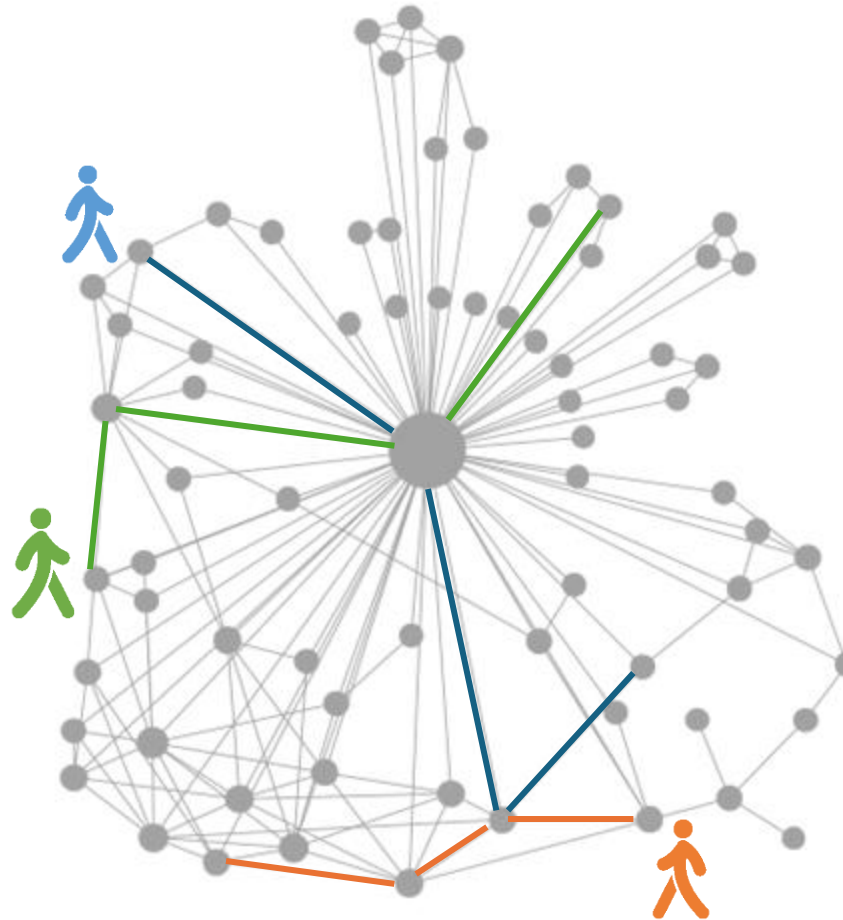| Location | L1 cache | L2 cache | L3 cache | Local mem | Remote mem |
|---|---|---|---|---|---|
| Sequential read | 0.42ns | 0.41ns | 0.44ns | 0.76ns | 1.51ns |
| Random Read | 0.77ns | 0.95ns | 2.60ns | 18.35ns | 24.35ns |
| Pointer-chasing | 1.69ns | 5.26ns | 19.26ns | 116.90ns | 194.26ns |

**4-byte read latency at different cache/DRAM layers**

☐ Sequential reads quite cheap, and relatively uniform
  ➢ Even across NUMA node (remote memory)

☐ Random reads slower, with wider distribution
  ➢ Large gap between L3 and DRAM
  ➢ Pointer-chasing especially costly: even in L3

- ## Problem definition
  - Input: graph, set of walkers placed at starting vertices
  - Each walker walks around
    - By randomly selecting an edge to follow
    - For given number of steps or till given termination condition
  - Output
    - Computation during walk, and/or
    - Set of walk paths
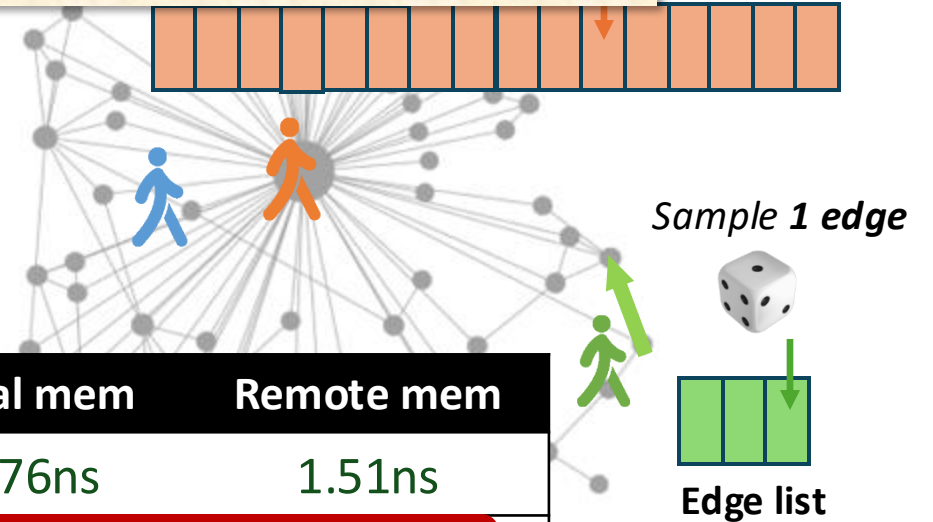
☐ Walke...

> ➤ F...

💡 Random walk doesn't mean **random memory accesses**

▪ Plenty of temporal and spatial locality around!

*...cheline discarded*

☐ Probl...

> ➤ Accept randomness in memory accesses
> ➤ Process walkers in turn, wherever they are
>   > ➤ High data dependency (pointer chasing)
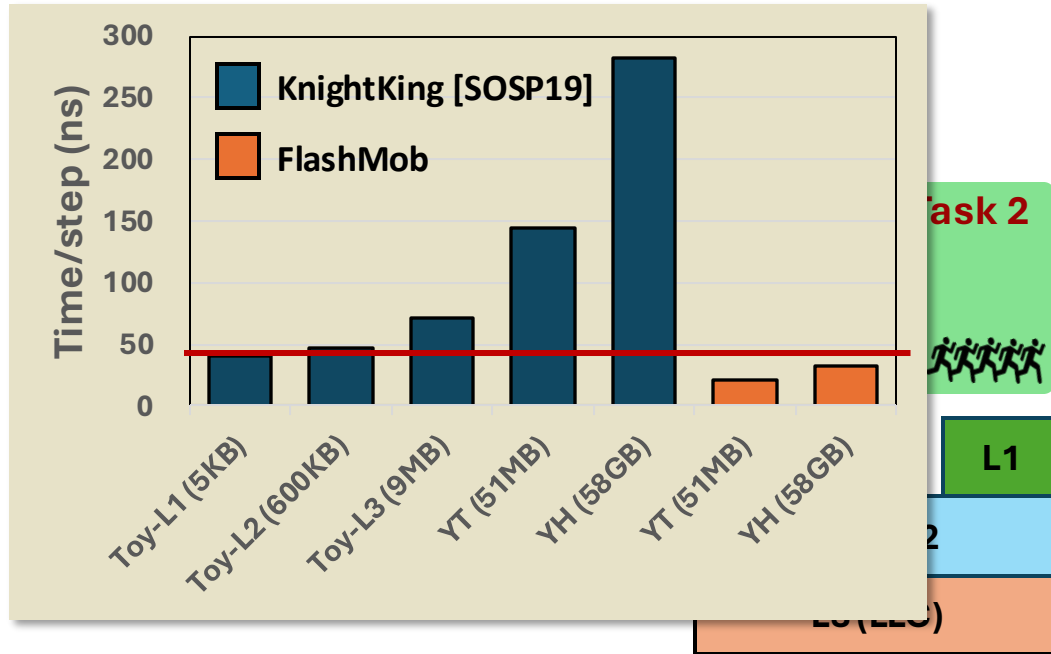>   > ➤ Low utilization of loaded cache lines

*Sample 1 edge*

**Edge list**

| Location | L1 cache | L2 cache | L3 cache | Local mem | Remote mem |
|---|---|---|---|---|---|
| Sequential read | 0.42ns | 0.41ns | 0.44ns | 0.76ns | 1.51ns |
| Random Read | 0.77ns | 0.95ns | 2.60ns | 18.35ns | 24.35ns |
| Pointer-chasing | 1.69ns | 5.26ns | 19.26ns | 116.90ns | 194.26ns |

**Access latency at different cache/DRAM layers**

# FlashMob [SOSP21]: Cache-speed Graph Random Walk



Bar chart — Time/step (ns) vs. graph/cache configurations:
- Legend: KnightKing [SOSP19] (dark blue), FlashMob (orange)
- X-axis: Toy-L1 (5KB), Toy-L2 (600KB), Toy-L3 (9MB), YT (51MB), YH (58GB), YT (51MB), YH (58GB)

**Task 2**

**L1**

**L3 (LLC)**

**Task 1**

**All walkers in partition**

**Cache-aware vertex partitions**

- High concurrency is good!
  - Economy of scale for objects cached

**DRAM**

❑ New challenges due to high-performance hardware

➢ Faster storage: 1000x to 10x latency gap from DRAM

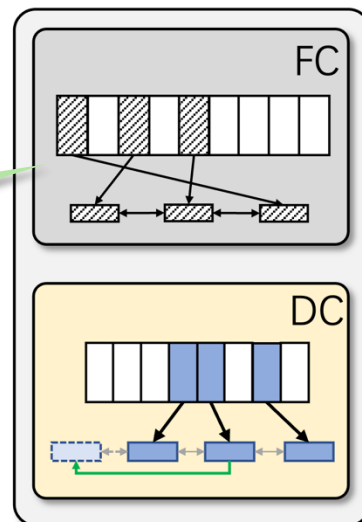➢ Scalability to high core counts

- High concurrency is bad!
  - ▪ Write contention on cache hits

❑ FrozenHot [EuroSys '23]

➢ Speeding up hit path by removing cache management

**Most hits on frozen cache (FC)**
- **No list management**
- **No locks**
- **Faster hash table**
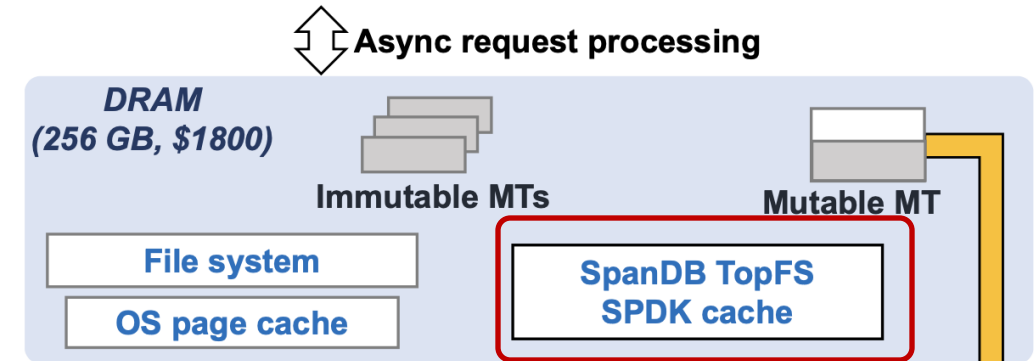
FC

DC

- ❏ SpanDB [FAST21]
  - ➢ Implemented within RocksDB

- ❏ Distributing LSM-tree based KV data
  - ➢ Large and slow disk for capacity
  - ➢ Small and fast disk for speed

- ❏ Automatic tree layer placement
  - ➢ Adaptive to partition sizes and workloads
  - ➢ Allow hybrid storage for cost-effectiveness



Async request processing

DRAM (256 GB, $1800)
Immutable MTs
Mutable MT
File system
OS page cache
SpanDB TopFS SPDK cache

NVMe SSD — 2TB, 500k IOPS

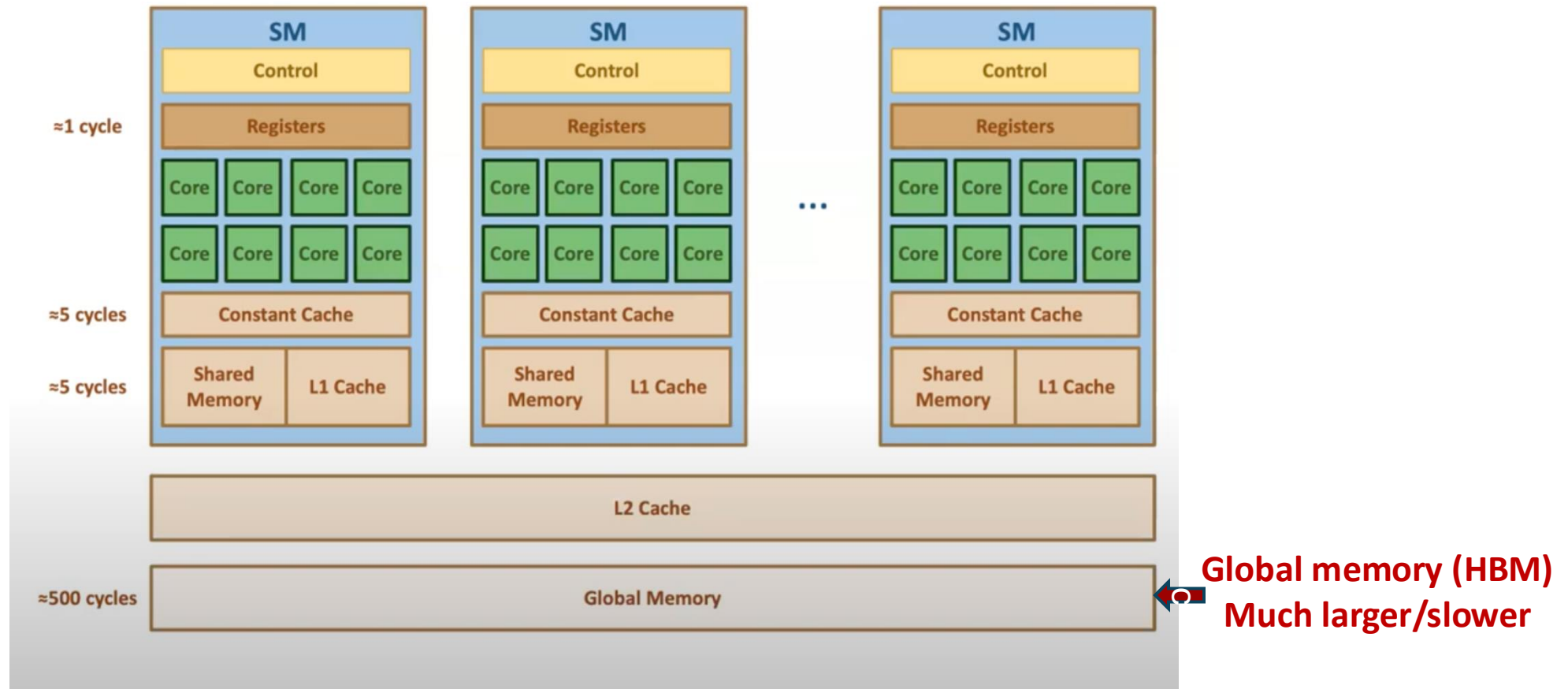NVMe SSD — 500GB, 1M IOPS, 2GB/s BW

SATA SSD — 5TB, 3GB/s BW

# Storage and I/O Not Efficiently Shared

❑ Major factor leading to storage frustrations and wastes

➢ From single server to data-centers/supercomputers

❑ In this talk

➢ Sample related problems and solutions in our past research

➢ Extended I/O hierarchy for AI workloads

# AI and LLM Age: GPU as a Supercomputer



**NVIDIA H100 Memory Hierarchy**

# GPU Global Memory the New Disk?

❑ FlashAttention: optimization targeting long sequences
  ➤ Adopted by major LLM frameworks: PyTorch, Megatron, DeepSeek ...

FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness

Tri Dao[†], Daniel Y. Fu[†], Stefano Ermon[†], Atri Rudra[‡], and Christopher Ré[†]

[†]Department of Computer Science, Stanford University
[‡]Department of Computer Science and Engineering, University at Buffalo, SUNY
{trid,danfu}@cs.stanford.edu, ermon@stanford.edu, atri@buffalo.edu, chrismre@cs.stanford.edu
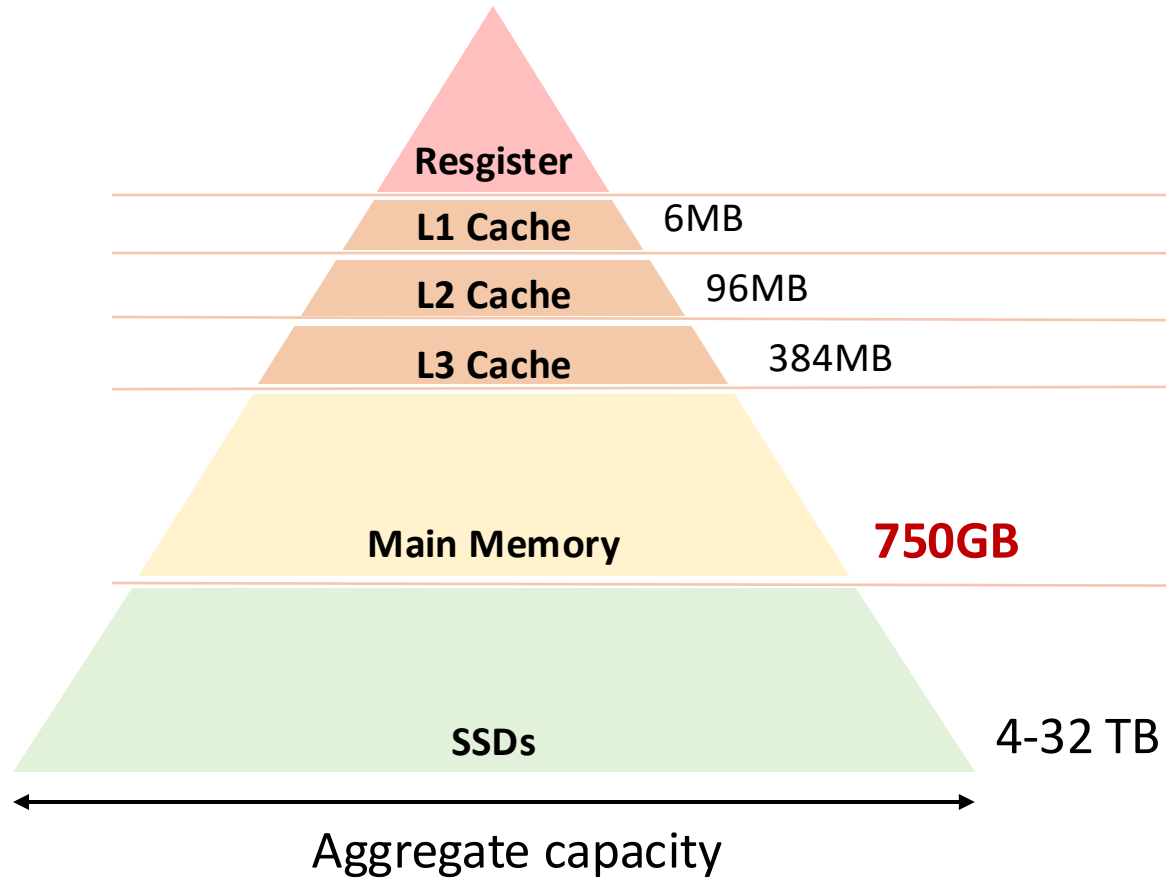
June 24, 2022

## Abstract

Transformers are slow and memory-hungry on long sequences, since the time and memory complexity of self-attention are quadratic in sequence length. Approximate attention methods have attempted to address this problem by trading off model quality to reduce the compute complexity, but often do not achieve wall-clock speedup. We argue that a missing principle is making attention algorithms *IO-aware*—accounting for reads and writes between levels of GPU memory. We propose FLASHATTENTION,

# CPU/GPU-Side Storage Hierarchy (Capacity)

**CPU (AMD EYPC 9654)**

- Resgister
- L1 Cache — 6MB
- L2 Cache — 96MB
- L3 Cache — 384MB
- Main Memory — **750GB**
- SSDs — 4-32 TB

Aggregate capacity

**GPU (H100)**

- Resgister
- L1 Cache — 33MB
- L2 Cache — 50MB
- High Bandwidth Memory — **80GB**

# CPU/GPU-Side Storage Hierarchy (Bandwidth)

**CPU (E.g. AMD EYPC 9654)**

Resgister

L1 Cache — 44TB/s

L2 Cache — 12TB/s

L3 Cache — 768GB/s

Main Memory — **460GB/s**

SSDs

← Aggregate bandwidth →

**GPU (E.g. H100)**

Resgister

L1 Cache — 33TB/s

L2 Cache — 12TB/s

HBM — **3.35TB/s**

# "I/O Characteristics" of Transformer Components

| Operator | Workload | Access type | Bottleneck |
|---|---|---|---|
| **Normalization** | Training & inference | Balanced R-W | Memory |
| **GEMM** | Training & inference | Mainly reads | Compute → Memory |
| **Attention** | Training & inference | Mainly reads | See table below |
| **Dropout** | Training | Balanced R-W | Memory |
| **Activation Function** | Training & inference | Balanced R-W | Compute → Memory |

| Attn sequence length | Workload | Bound |
|---|---|---|
| Long | Training & prefill | Compute → Memory |
| Short | Training & prefill | Compute → Memory |
| **Long** | **Decode** | **Memory** |
| Short | Decode | Compute → Memory |

# Good News: Storage-Friendly Access Patterns

❑ Large, sequential, read-heavy accesses
  ➢ Little random reads, relatively light writes
  ➢ Regular, predictable, collaborative data streaming
  ➢ **Data content/precision could be manipulated!**

❑ Many storage/HPC tricks apply
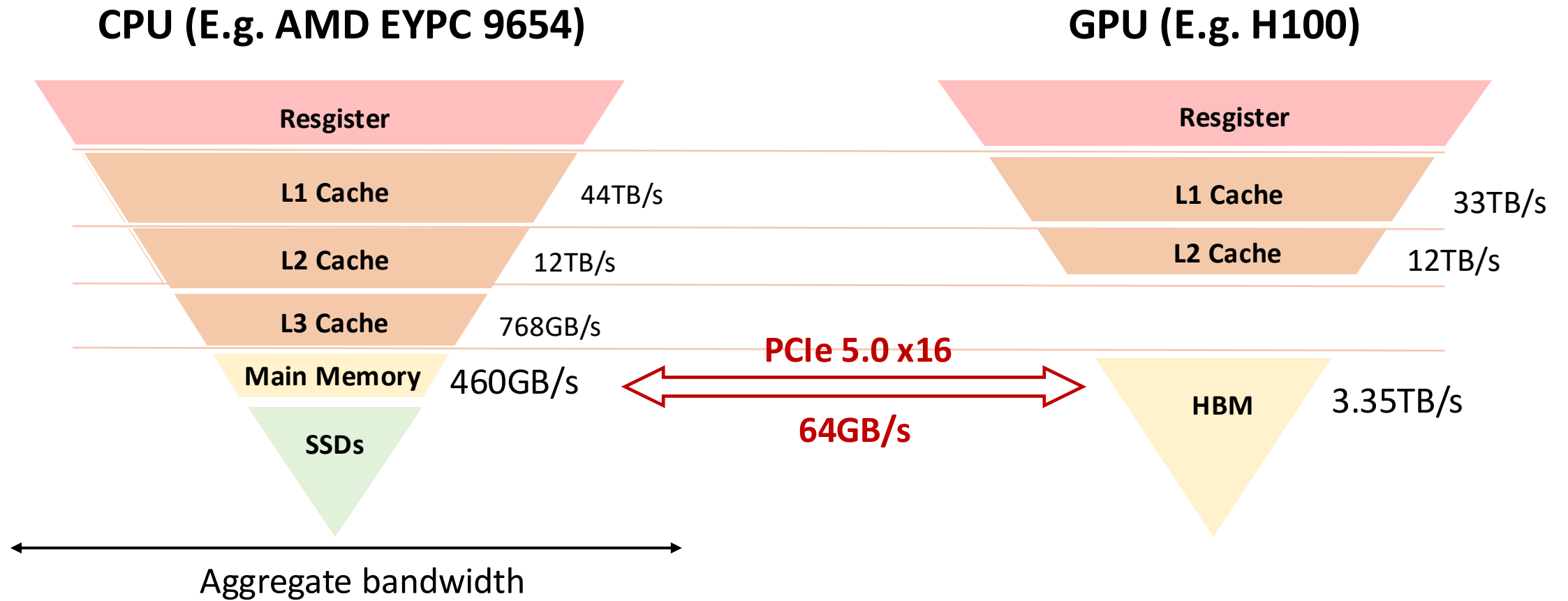  ➢ Prefetching
  ➢ Tiling/Tiering

# **Bad News: High Efficiency Demands Large Space**

❑ GPU MFU (Model FLOPS Utilization) relies on batch size

➤ Larger batches -> higher parallelism, more data reuse

➤ Current leading frameworks get <50% of GPU peak TFLOPS

❑ Batch size limited by HBM size

➤ Especially w. **long-sequence attention** in **decoding**

# Can CPU-Side Memory Help?

**CPU (E.g. AMD EYPC 9654)**

**GPU (E.g. H100)**

Resgister

L1 Cache     44TB/s

L2 Cache     12TB/s

L3 Cache     768GB/s

Main Memory     460GB/s

SSDs

Resgister

L1 Cache     33TB/s

L2 Cache     12TB/s

HBM     3.35TB/s

**PCIe 5.0 x16**

**64GB/s**

Aggregate bandwidth

# Recent Work on Reducing KV-Cache I/O Demands

❏ Parameter/KV Cache offloading
  ➢ FlexGen [ICML23]
  ➢ MoE-Lightning [ASPLOS25]

❏ KV Cache Compression
  ➢ Keyformer [MLsys24]

❏ Quantization
  ➢ ZipCache [NeurlPS24]

❏ Window attention
  ➢ StreamingLLM [ICLR24]

# Closing Remarks

❑ **Shared nature makes storage challenging and interesting**

  ➢ Contention and interference, but also higher throughput and utilization

  ➢ Joint CPU-GPU storage hierarchy creates more scenarios for sharing/coordination

    • HBM too small to saturate GPU cores, too fast for DRAM to stream

❑ **Education also challenged by new modes of knowledge sharing**

  ➢ *AI practitioners need to know systems basics*

  ➢ *CS students need to retain focus/courage in system building*