## Storage Class Memory is Dead, All Hail **M**anaged-**R**etention **M**emory: Rethinking Memory for the Al Era

Sergey Legtchenko, Ioan Stefanovici, Richard Black, Ant Rowstron, Junyi Liu, Paolo Costa, Burcu Canakci, Dushyanth Narayanan, Xingbo Wu

Microsoft Research

HotOS'25

#### Al Inference – The Dominant Cloud Workload

Generative AI has changed the game & inference demand is huge

**NEWS** 10 April 2025

## Data centres will use twice as much energy by 2030 — driven by AI

These facilities accounted for roughly 1.5% of global electricity consumption in 2024.

By Sophia Chen









# Tech megacaps plan to spend more than \$300 billion in 2025 as Al race intensifies

PUBLISHED SAT. FEB 8 2025-8:00 AM EST | UPDATED SAT. FEB 8 2025-11:02 AM EST





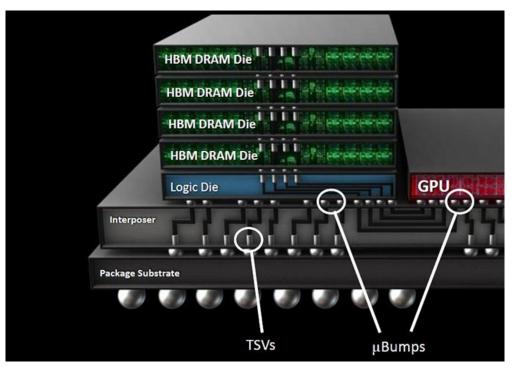
A great challenge and opportunity for the systems community to rethink systems architecture

#### The Problem Today: The Curse of HBM

**H**igh **B**andwidth **M**emory is the **only** option today to achieve good bandwidth to Al data

**But**, a litany of problems...

- Complex manufacturing and packaging
- Unreliable
- Expensive: significant portion of GPU cost and power



source: https://www.anandtech.com/show/9969/jedec-publishes-hbm2-specification

#### What is HBM in *LLM Inference* Actually Used For?

LLM inference: the prominent workload

Two large data structures: model weights + KV cache

Model weights (~ 50% today) KV-cache (~ 50% today):

Write: once Write: append-only

Read: each forward pass Read: in whole each forward pass

Observation: very large, predictable, sequential Reads dominate

- HBM is "overprovisioned" on write performance
- Random access of HBM not necessary

Can we leverage the specific properties of Al inference to design a better memory?

#### A New Class of Memory for Al Inference

- "New" memory technologies: STT-MRAM, ReRAM, PCM, FeRAM,...
  - Viewed through "Storage Class Memory" lens long-term data retention was a goal
- For Al Inference:

Important Metrics	Less Important Metrics
Capacity / \$	Write performance
Read bandwidth	Random access
Energy	Long-term Retention

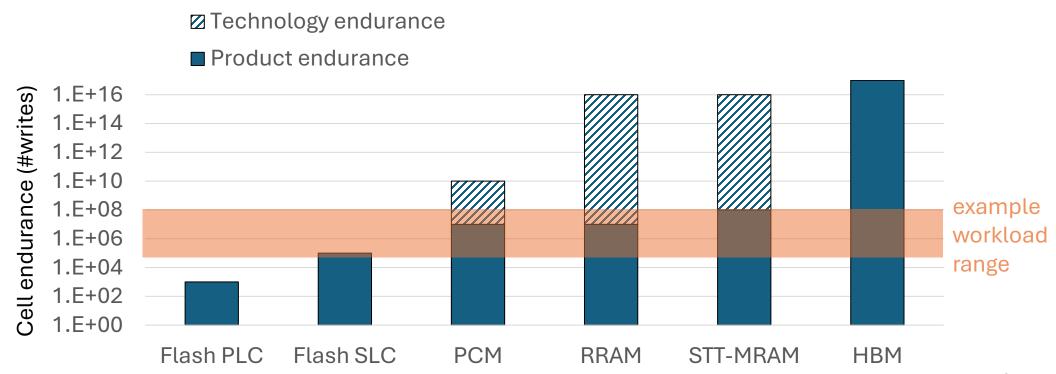
**Key insight:** possible to trade-off *write performance & retention time* for important metrics

- Storage Class Memory non-volatility (10+ yr retention) is not required
- Hours-long retention time is sufficient and enables power advantage

Managed-Retention Memory: a new class of memory for Al inference

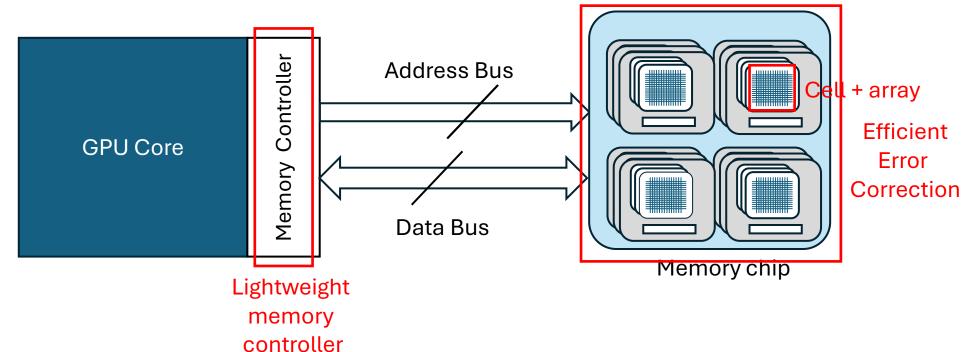
### MRM: A New Opportunity for SCM Technologies

Existing memory technologies **can** inherently be optimised for MRM example trade-off: retention \( \gamma \) endurance \( \psi



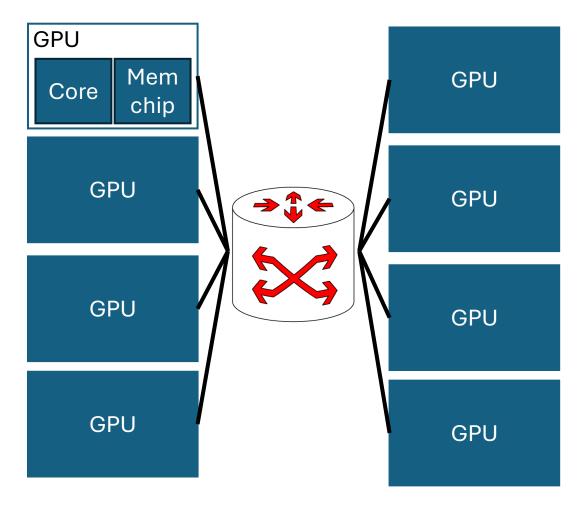
#### MRM Research Opportunities

- Innovation across the HW/SW stack needed
- How do we take leverage the workload?
  - lack of random access, lack of refreshes



#### MRM Systems Opportunities

- MRM abstraction: how to expose MRM to systems?
- Dynamically configurable retention
  - Should software configure retention period per write?
- Retention-aware data placement & scheduling
  - Software-driven movement



#### MRM: Rethinking Memory for the Al Era

1) Massive opportunity and need to disrupt HBM for Al inference

- 2) Managed-Retention Memory: a new class of memory that trades off retention and write performance for energy, read performance, and cost
- 3) Ripe for innovation across cells, arrays, controllers, system abstractions, and much more!