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

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

Microsoft Research

# **AI Inference** – The Dominant Cloud Workload

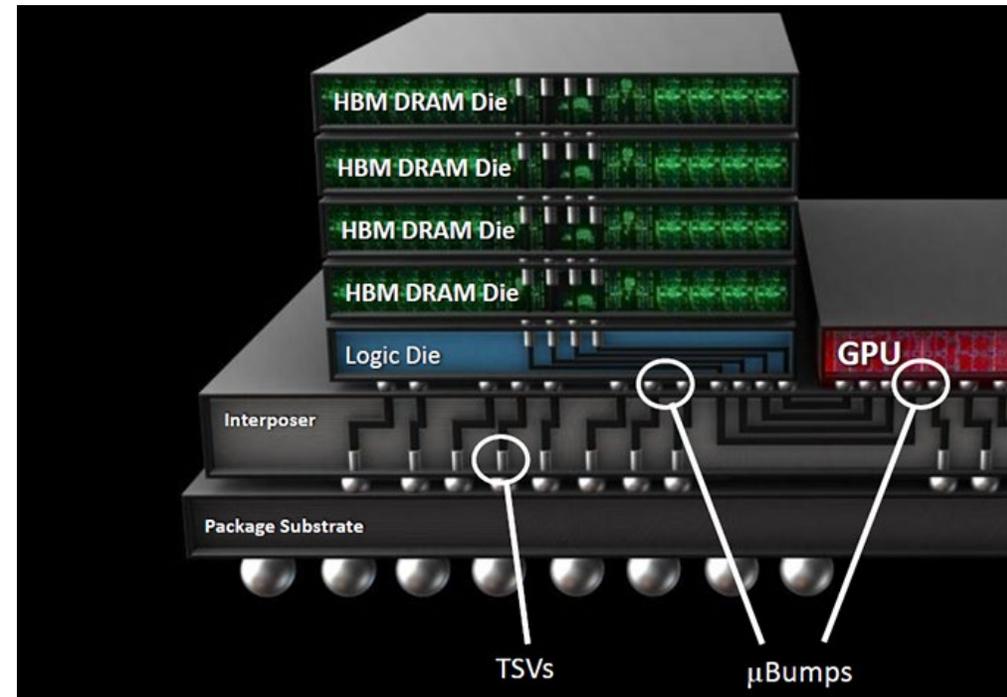Generative AI has changed the game & inference demand is huge

A great challenge and opportunity for the research community to rethink system hardware and software architectures

# The Problem Today: The Curse of HBM

**H**igh **B**andwidth **M**emory is the **only** option today to achieve good bandwidth to AI 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)
- Write: once
- Read: each forward pass

KV-cache (~ 50% today):
- Write: append-only
- Read: in whole each forward pass

Observation: very large, predictable, block Reads dominate
- HBM is "overprovisioned" on write performance
- Small, random access of HBM not necessary

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

# A New Class of Memory for AI Inference

- "New" memory technologies: STT-MRAM, ReRAM, PCM, FeRAM,…
    - Viewed through "Storage Class Memory" lens – long-term data *retention* was a goal

- For AI Inference:

| Important Metrics | Less Important Metrics |
|---|---|
| Capacity / $ | Write performance |
| Read bandwidth | Small, random access |
| Energy | **Long-term Retention** |

**Key insight:** possible to trade-off *write performance & retention time* for important metrics
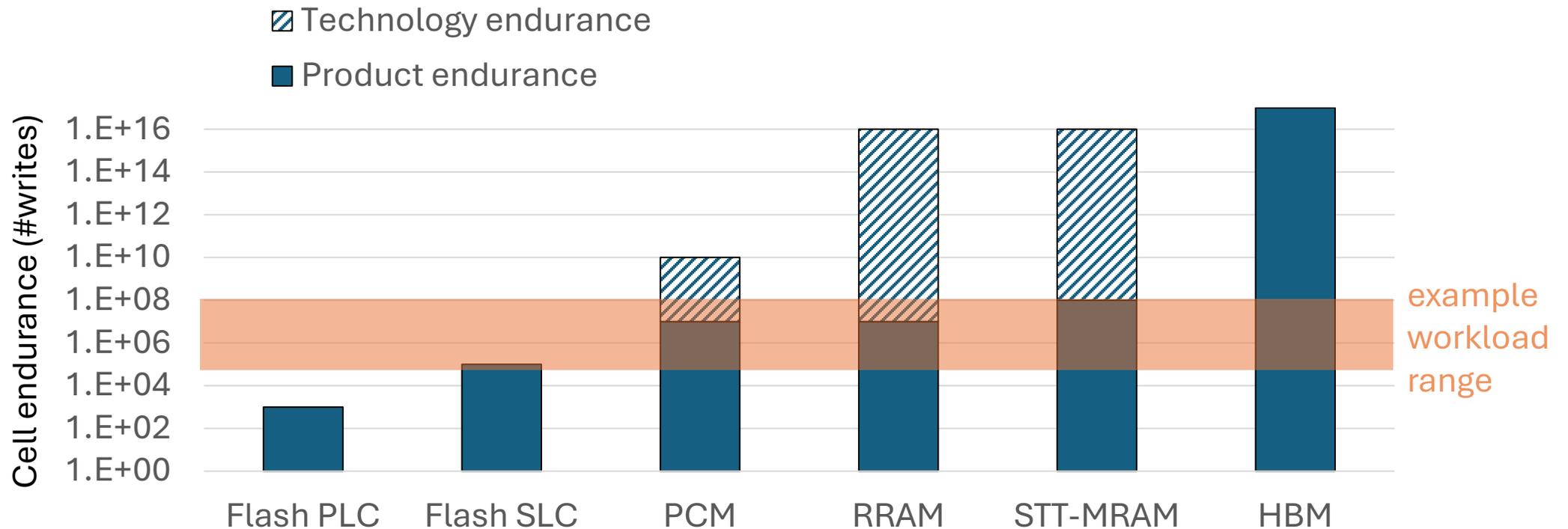- Storage Class Memory non-volatility (10+ yr retention) is not required
- Hours-long retention time is sufficient and enables power advantage

**M**anaged-**R**etention **M**emory: a new class of memory for AI inference

# MRM: A New Opportunity for SCM Technologies
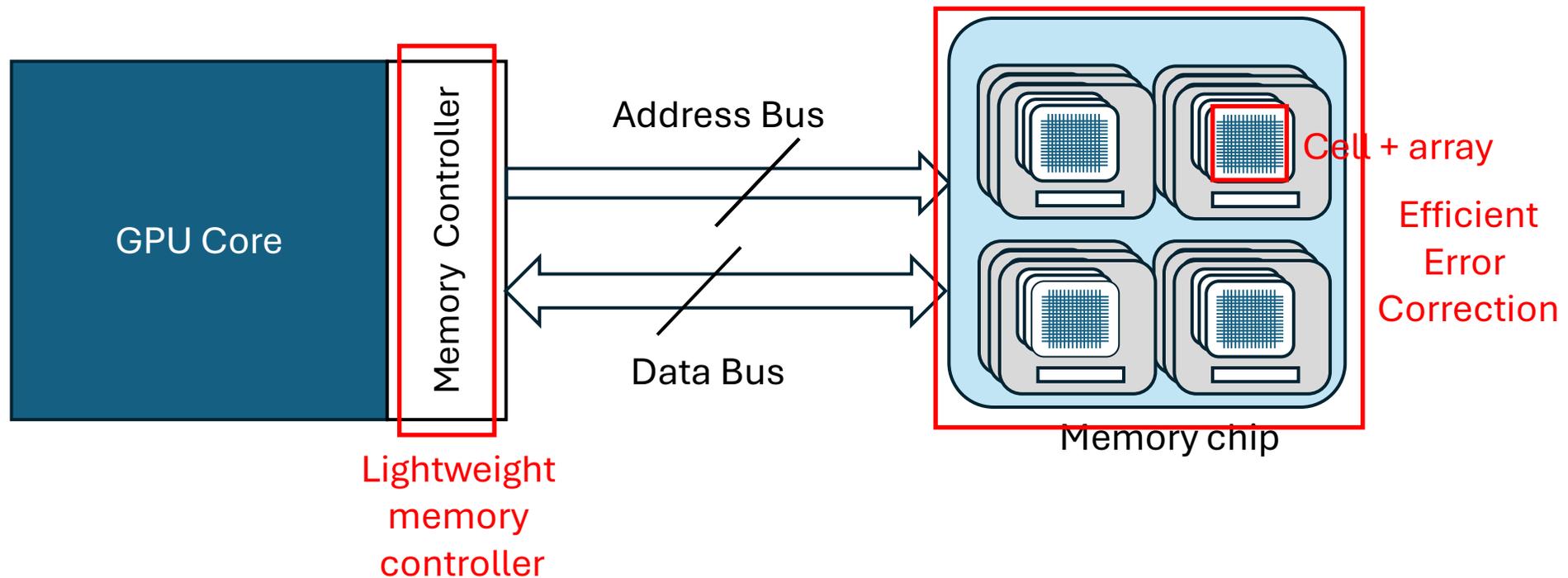
Existing memory technologies **can** inherently be optimised for MRM

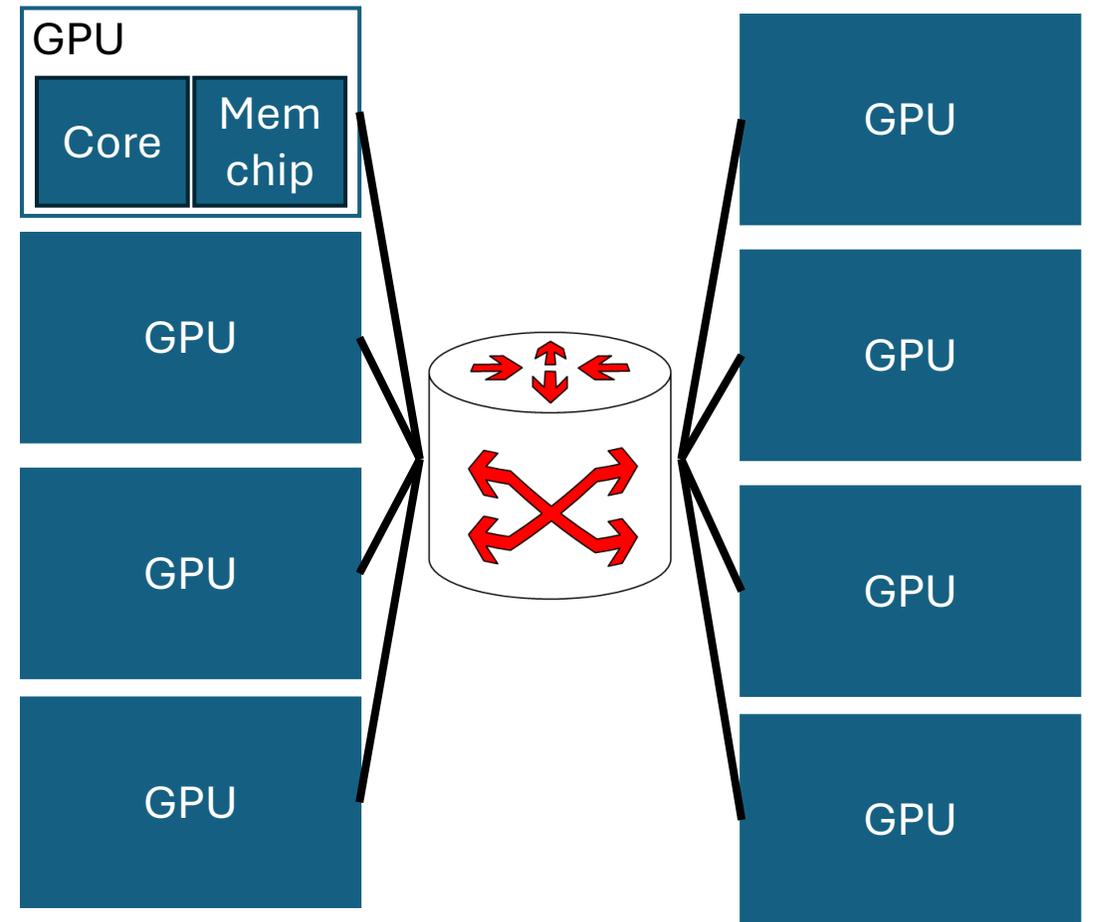example trade-off:        retention ↑    endurance ↓

# MRM Research Opportunities

- Innovation across the HW/SW stack needed

- How do we leverage workload + cell properties?
  - lack of small random access, lack of refreshes



GPU Core

Memory Controller

Address Bus

Data Bus

Lightweight memory controller

Cell + array

Efficient Error Correction

Memory chip

# MRM Systems Research 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 AI Era

1) <u>Massive</u> opportunity and need to disrupt HBM for AI inference

2) **M**anaged-**R**etention **M**emory: 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!