

# Metadata Management for Large-Scale Hybrid Memory Architectures

## Introduction

### Requirements for main memory

- More bandwidth and shorter latency for high performance applications
  - More memory capacity for big data
- Both high performance and a large capacity are required for main memory.

### Hybrid memory architectures (HMAs)

Combine two types of memories to realize both high performance and large capacity

#### Far memory (FM)

Large capacity but low performance to store infrequently-used data

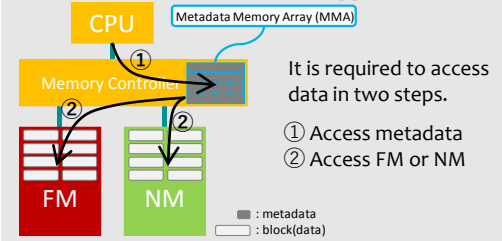
#### Near memory (NM)

High performance but small capacity to store frequently-used data

Management of data between NM and FM is essential to keep many frequently-used data in NM.

### Metadata

Records the management of information such as location of the corresponding data [1]

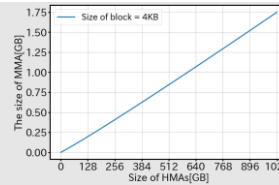


## Problem : Increasing the Capacity for Metadata on an HMA

Estimation of the size of metadata memory array (MMA) in the future

$$\text{(The size of MMA)} = \frac{(\text{Size of HMAs})}{(\text{Size of block})} \times \left( 2 \times \frac{1}{8} \log_2 \left( \frac{(\text{Size of HMAs})}{(\text{Size of block})} \right) \right) [\text{Bytes}]$$

The number of metadata      Size of metadata



- The size of MMA increases as the size of the HMA increases.

It is not practical to store metadata only on a chip in the future.

## Analysis of Metadata Access Pattern

### In-Memory metadata with metadata memory array

#### Idea

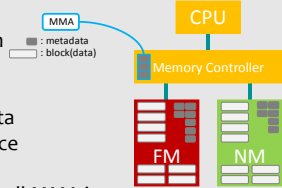
- Manage metadata by distributing them to FM/NM and memory controller

#### Problem

- The long access latency to the metadata in FM/NM may degrade the performance

#### Solution

- Stores only necessary metadata in a small MMA in a memory controller to avoid performance degradation



### Behavior of metadata accesses

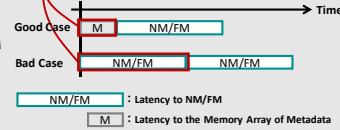
Equation of the time to access data

$$\text{(Time to access data)} = \text{(Time to refer metadata)} + \text{(Time to access HMAs)}$$

Good case : Access metadata in MMA

Bad case : Access metadata in NM/FM

Bad case increases the access latency to the data.



To find a way to save metadata in the MMA in advance, it is required to analyze metadata access patterns.

## Experimental Results and Discussions

### Purpose

To analyze block access patterns for efficient control of metadata

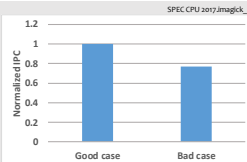
### Experimental environment

Obtain an address trace based on block address using simulators (Gem5 and NVMain) and analyze the accesses to metadata

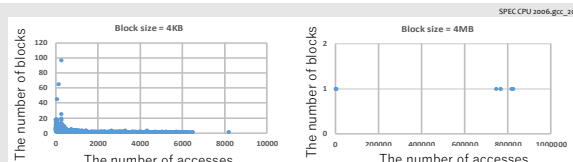
- [Benchmark] SPEC CPU 2006, 2017
- [Cache size] L1:64KB, L1D:64KB, L2:512KB, L3:2MB

- [Latency] M:2cycles, NM:32cycles, FM:68cycles
- If data is stored in NM/FM, NM/FM is accessed to refer metadata.

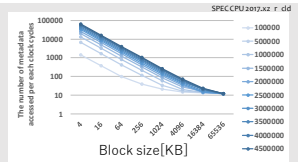
### Results



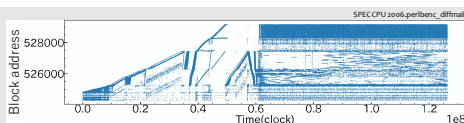
- The bad case decreases IPC by accessing metadata in NM and/or FM many times.
- It is important for performance to predict necessary metadata and store the predicted metadata in the MMA in advance.



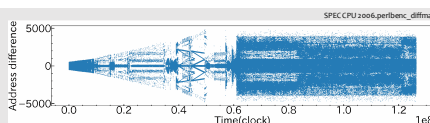
- When the block size is large, the number of blocks is limited.
- By using the large block size, blocks can be categorized by the number of accesses.



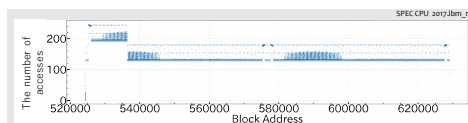
- By decreasing clock cycles and/or increasing the size of block, the number of metadata decreases.



- Frequently accessed data or its neighbors continue to be frequently accessed.
- Infrequently accessed data will not be accessed.



- The address between currently accessed data and next accessed data is within a certain range.



- Neighborhood of data have the similar number of accesses.
- By using the block address, it is possible to predict the number of accesses.

## Conclusions & Future Work

- Large-scale HMAs need more metadata, and the capacity for metadata should be mitigated.
- This research analyzes block access pattern to realize the in-memory metadata management with metadata caching.
- The experimental results show the behavior of metadata accesses to predict metadata in the HMAs.
- In the future, it is planned to implement the metadata prediction mechanism by using the metrics.

### Reference

[1] Chia Chen Chou, et al. Cameo: A two-level memory organization with capacity of main memory and flexibility of hardware-managed cache. In Proc. of Micro, pp. 1-12, 2014.