ROSS:

A Design of Read-Oriented STT-MRAM Storage for Energy-Efficient Non-Uniform Cache Architecture

Jie Zhang, Miryeong Kwon, Changyoung Park, Myoungsoo Jung, Songkuk Kim Computer Architecture and Memory System Lab, Yonsei University





Summary

Motivation:

There is a huge demand to increase the last-level cache size. However, SRAM suffers from the **high leakage power** and **low density**.

STT-MRAM can be a good candidate to replace SRAM, due to its small cell size and low leakage power.

Challenge:

Conventional STT-MRAM cannot directly substitute SRAM, because of

- High write latency

- High write energy

Solution:

We observe read-oriented data pattern in real applications; We propose **ROSS**, a read-oriented STT-MRAM storage that can

- dynamically detect the read-oriented data pattern
- reduce the write frequency in STT-MRAM
- accommodate as many read misses as possible in LLC

Motivation

Data set size increases explosively.



Motivation

Total Cache Size per Chip (L2 + L3) vs. Time



Huge demand to increase on-chip last-level cache size in future

Challenges

Traditional SRAM



Challenges

	Traditional SRAM	STT-MRAM	
		SL MTJ Low-Resistance High-Resista State Free Layer Oxide Fixed Layer Parallel Configuration Configuration	L ince r el on
Read Latency	397ps	<u>238ps</u>	
Write Latency	397ps	6932ps	High Write
Read Energy	35pJ	<u> </u>	latency/Energy
Write Energy	35pJ	90pJ	
Cell Area (F^2)	50~120	6~40	High density
Leakage	75.7mW	<u>6.6mW</u>	Low leakage

Prior work

<u>Read Preemption</u>: suspend write operation and give priority to read operation

- **D** Mitigate the impact of long write latency to read access;
- Cannot reduce the write latency

Write Request Rerouting: accommodate write requests with idle blocks

- Parallelize write requests to hide long write latency;
- □ Cannot reduce the write energy;



ROSS-key observation

How the application behavior can impact LLC data access?



What if we can detect the singular write blocks and put in STT-MRAM?

sopx-calix-

gbmk-

namd-

cacs

gmcs

oer-

bzipame

milc

mcf

mer

sjng-

Jems-

-mq omn2

ROSS

- Target: build a hybrid last-level cache to
- i. avoid many cache read misses with larger STT-MRAM;
- ii. accommodate write updates in SRAM;
- iii. intelligent data migration algorithm;

Hybrid-NUCA



Hybrid-NUCA

Address Migration policy map policy LLC Controller SRAM-SRAM migrat. arrays 1-way **SRAM** banks Cache SRAM->STT-MRAM STT-MRAM banks -way migration (M) **Cache Bank Set**

Less flexibility, but simplified design

Hybrid-NUCA



Request routes to specific cache bank set

Tag broadcast to each bank for searching

Data Flow



Singular write detection







Initialize

Poll

Evict

Experiment Setup

Evaluation model:

ND-NUCA: traditional SRAM NUCA model;

DB-NUCA: STT-MRAM NUCA model aware of deadblock;

HB-ROSS: Hybrid Read-Oriented STT-MRAM Storage;

LER-ROSS: Early Retirement aware Raed-Oriented STT-MRAM Storage;

Gem5 simulation parameters and workloads:

Processor	1 core, OoO execution, SE mode	
Frequency	Frequency 3.2GHz	
L1 Cache	16KB/core, 2-way, 2 cycles	
SRAM L2	32KB/bank, R/W : 20 cycles	
STT-MRAM L2	96KB/bank, R/W : 20/60 cycles	
Network	Wormhole Switching, 2 cycles	
Off-chip Mem	ff-chip Mem 4GB DDR3 DRAM FR-FCFS	
Benchmark	SPEC2006	
Workloads	PerlBench, mcf, milc, libquantum, lbm,	
WOLKIOAUS	cactusADM, sjeng, and gobmk	

Evaluation



DB-NUCA: suffer from long write latency

HB-ROSS: larger capacity and less STT-MRAM write

ER-ROSS: benefits from early retirement

Evaluation



DB-NUCA: suffer from evicting singular write blocks

ROSS: STT-MRAM accommodates read-intensive data

Evaluation



ND-NUCA: consume more leakage energy

ROSS: reduce the write energy overhead in STT-MRAM

Thank You