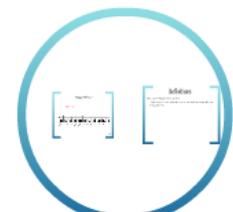
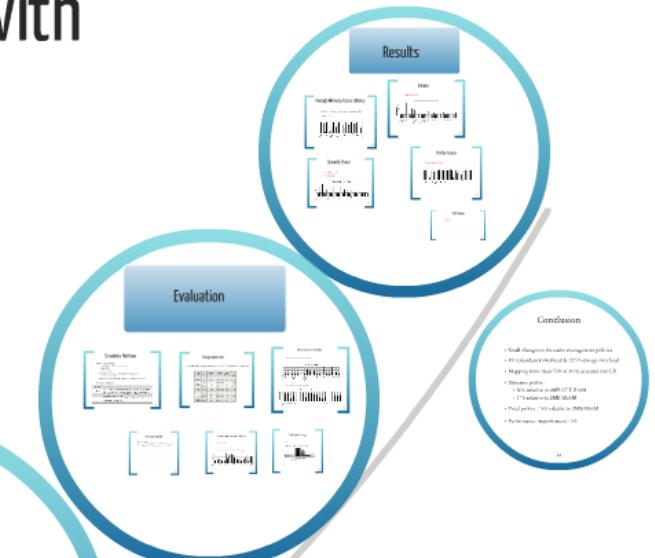
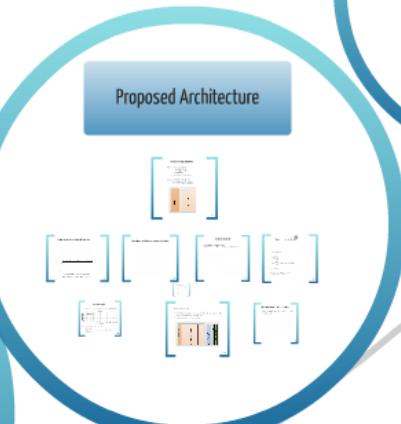
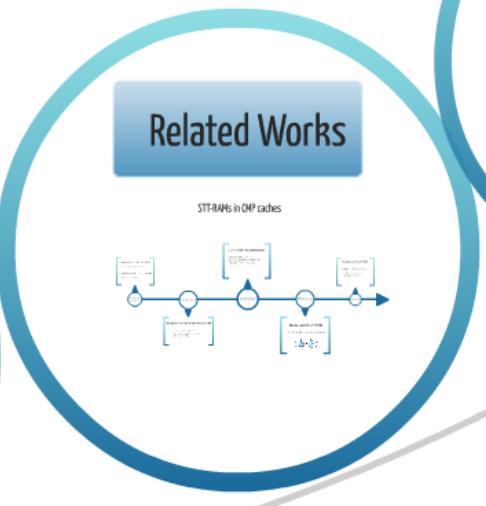
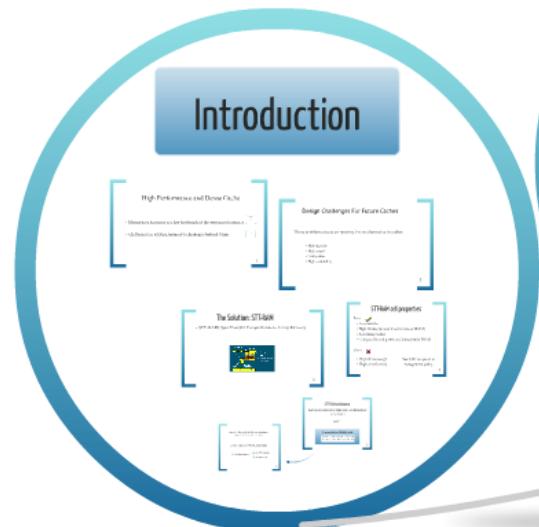


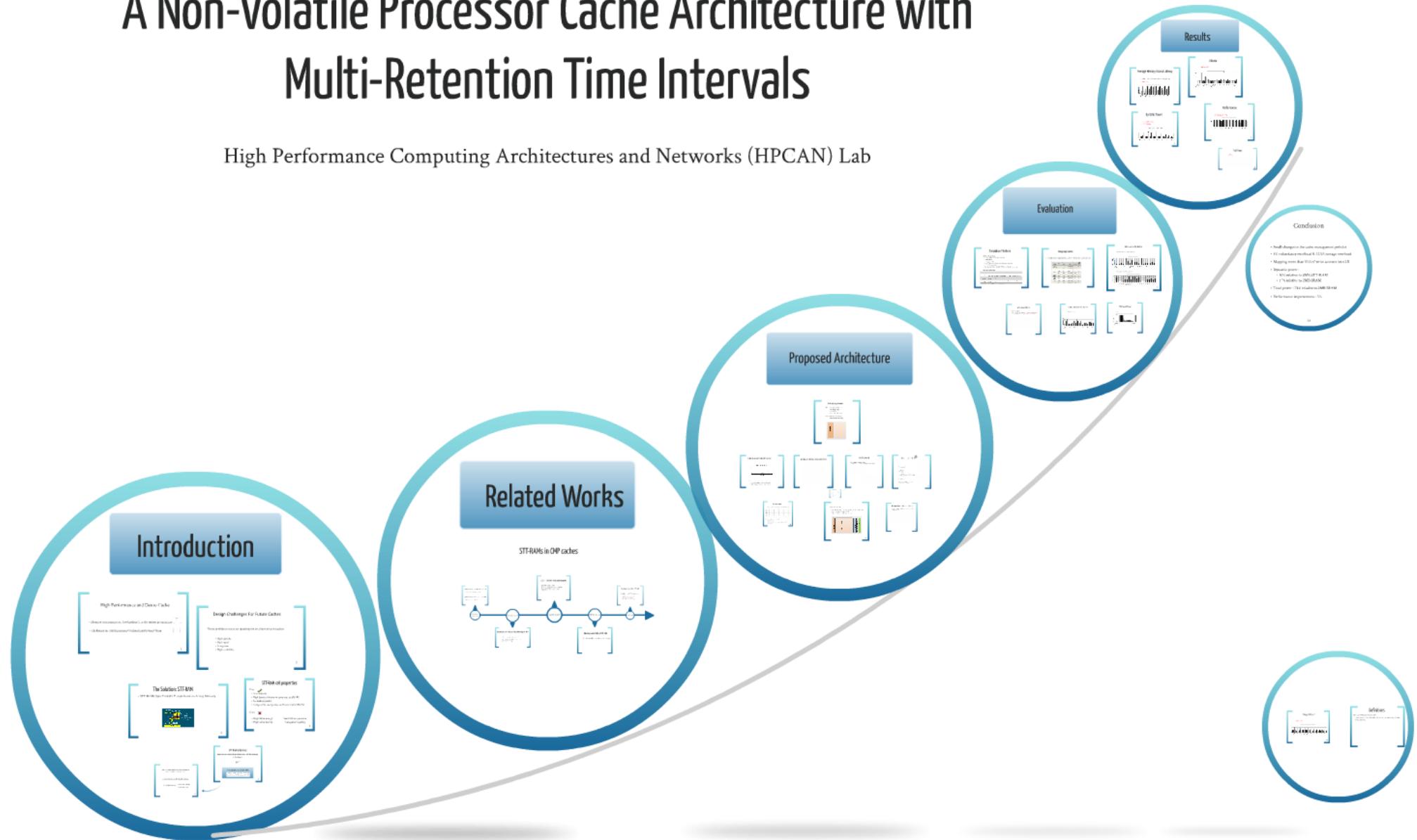
A Non-Volatile Processor Cache Architecture with Multi-Retention Time Intervals

High Performance Computing Architectures and Networks (HPCAN) Lab



A Non-Volatile Processor Cache Architecture with Multi-Retention Time Intervals

High Performance Computing Architectures and Networks (HPCAN) Lab



Introduction

High Performance and Dense Cache

- Memory performance as a key bottleneck of the system performance
- Challenges for existing memory technologies beyond 45nm

1

Design Challenges For Future Caches

These problems create an opening for an alternative in caches

- High-density
- High speed
- Low power
- High-scalability

2

The Solution: STT-RAM

- STT-RAM: Spin Transfer Torque Random Access Memory



3

STT-RAM cell properties

- Pros ✓
- Non-volatile
 - High density (4x size in same area as SRAM)
 - No leakage power
 - Comparable read power and latency with SRAM

Cons ✗

- High write energy
- High write latency

Needs write operation management policy

4

STT-RAM endurance

much more better than other non-volatile memory techniques

BUT

Not perfect as SRAM cells

10^9 for DRAM vs 10^7 for SRAM - about 100 times worse for DRAM under 10^9

5

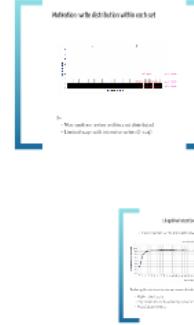
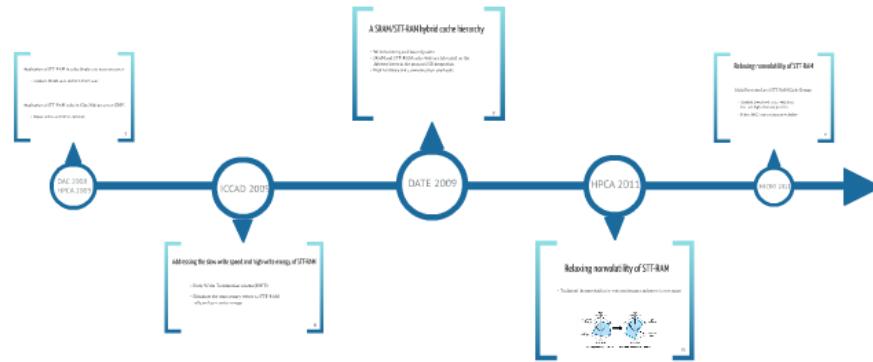
6

Using STT-RAM as LLC or SPC could benefit from write-back buffering by LLC cache

LL cache can use a direct-mapped scheme
• 2 Way set associativity
• LL write-through
• LL write-back

Related Works

STT-RAMs in CMP caches



Proposed Architecture



Ideal cache requirements:

- High utilization (75% to 80% usage)
- Low miss rate
- Low read time and energy

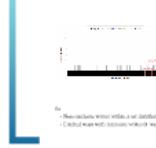
Low utilization (75% to 80% usage)

- Low access time and energy



13

Mutation with distribution within each set



14

Mutation with distribution between different sets



15

Best time interval



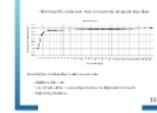
16

Always empty of 8 blocks



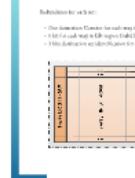
17

Updated writer file



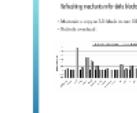
18

Reordering for S1 set



19

Reordering each byte for data blocks in S1 register



20

Evaluation

Simulation Platform

System-level simulation

- Virtutech SIMOS full-system simulator
- 4 cores
- Cycle accurate
- + CACTI timing model rounded by system latency
- + CACTI energy model
- Benchmark: ROI of the SPEC CPU2006 + PARSEC-2 programs

Simulator configuration:

Processor	Intel Xeon E5-2620 v2	Processor clock	2.1 GHz
L1 Cache	32 KB per core, 8 MB total, 4-way associativity, 128B cache line size, 10 ns access time	Latency	10 ns
L2 Cache	128 KB shared, 8-way associativity, 128B cache line size, 10 ns access time	Latency	10 ns
L3 Cache	16 MB shared, 4-way associativity, 128B cache line size, 10 ns access time	Latency	10 ns
Memory	4 GB DDR3-1600, 16-bit wide memory bus	Latency	10 ns

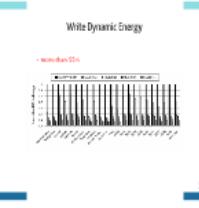
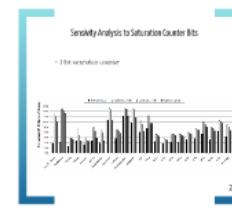
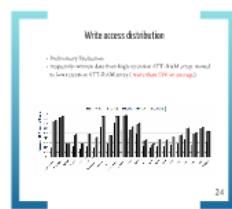
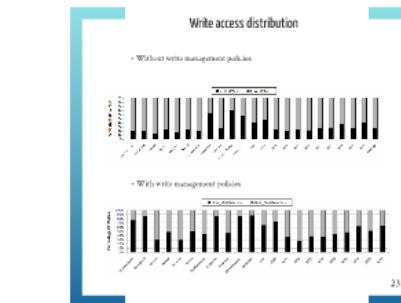
21

Design approaches

+ 1.3 2MB SRAM replacement with 8MB STT-RAM with same area

Cache Configuration	Op.	Execution Time	Dramatic Energy	Power	Area
SRAM 1.3M	R	1.000 cycles	0.000 J	0.000 W	0.000 cm ²
SRAM 8MB	R	—	0.000 J	—	0.000 cm ²
SRAM Modeling	R, W	1.000 cycles	0.000 J	0.000 W	0.000 cm ²
SRAM Modeling	R, W	1.000 cycles	0.000 J	0.000 W	0.000 cm ²
SRAM Modeling	R, W	1.000 cycles	0.000 J	0.000 W	0.000 cm ²
SRAM Modeling	R, W	1.000 cycles	0.000 J	0.000 W	0.000 cm ²
SRAM Modeling	R, W	1.000 cycles	0.000 J	0.000 W	0.000 cm ²
SRAM Modeling	R, W	1.000 cycles	0.000 J	0.000 W	0.000 cm ²

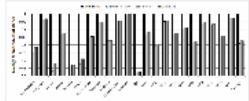
22



Results

Average Memory Access Latency

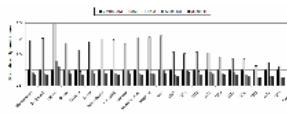
- AML = hit time + miss rate x miss penalty
 - AML {20%}



28

Dynamic Power

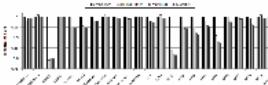
- 50% (STT-RAM)
 - 15% (SRAM)



1

Performance

- Performance (5%)



Total Power

-

Introduction

High Performance and Dense Cache

- Memory performance as a key bottleneck of the system performance
- Challenges for existing memory technologies beyond 45nm

1

Design Challenges For Future Caches

These problems create an opening for an alternative in caches

- High-density
- High speed
- Low power
- High-scalability

2

The Solution: STT-RAM

- STT-RAM: Spin Transfer Torque Random Access Memory



3

STT-RAM cell properties

- Pros ✓
- Non-volatile
 - High density (4x size in same area as SRAM)
 - No leakage power
 - Comparable read power and latency with SRAM

Cons ✗

- High write energy
- High write latency

Needs write operation management policy

4

STT-RAM endurance

much more better than other non-volatile memory techniques

BUT

Not perfect as SRAM cells
 10^9 for DRAM vs 10^7 for SRAM
become even for DRAM under 10^6

5

6

Using STT-RAM as LLC or SPC could benefit from write-back buffering by LLC cache.
LL cache can use a direct-mapped
• 2 Way set associativity
• 11 bit set address
• 11 bit word address

Introduction

High Performance and Dense Cache

- Memory performance as a key bottleneck of the system performance
- Challenges for existing memory technologies beyond 45nm



Performance Bottleneck

- Ever increasing number of on-chip cores
- Deep pipelining in microprocessor architectures
- Ever-increasing trend of memory footprint of the programs
(980MB for SPEC CPU2k6 vs. up to 200MB for SPEC CPU2k)
- Critical applications are becoming more data-centric,
less compute-centric

High Performance and Dense Cache

- Memory performance as a key bottleneck of the system performance
- Challenges for existing memory technologies beyond 45nm



Challenges beyond 45 nm

SRAM:

- High power consumption
- Leakage increasing 10X with each technology node
- Low scalability

High Performance and Dense Cache

- Memory performance as a key bottleneck of the system performance
- Challenges for existing memory technologies beyond 45nm



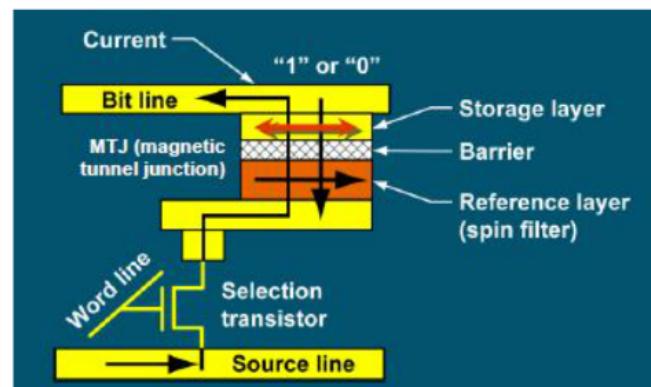
Design Challenges For Future Caches

These problems create an opening for an alternative in caches

- High-density
- High speed
- Low power
- High- scalability

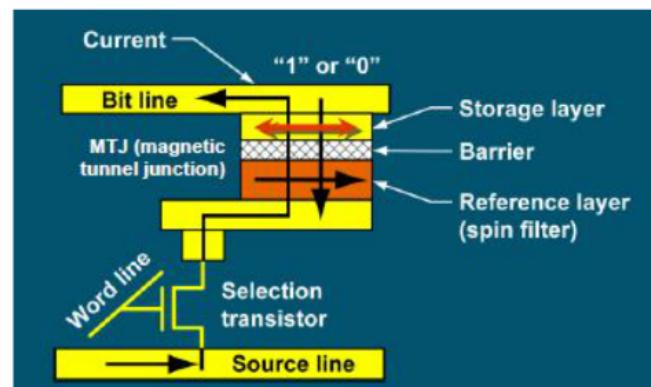
The Solution: STT-RAM

- STT-RAM: Spin Transfer Torque Random Access Memory



The Solution: STT-RAM

- STT-RAM: Spin Transfer Torque Random Access Memory



STT-RAM cell properties

Pros

- Non-volatile
- High density (4x size in same area as SRAM)
- No leakage power
- Comparable read power and latency with SRAM



Cons



- High write energy
- High write latency

Needs write operation
management policy

STT-RAM endurance

much more better than other non-volatile memory technologies

BUT

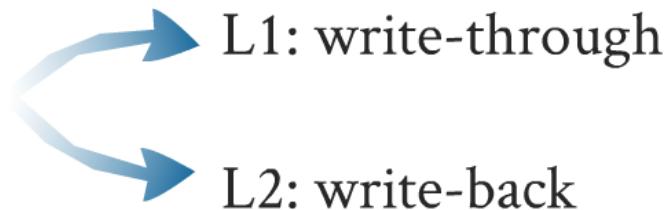
Not perfect as SRAM cells

*10^{16} for SRAM vs. 10^{15} for STT – RAM
become worse for MLC cells: under 10^{12}*

Using STT-RAM as LLC in CMP could benefit from write traffic filtering by L1/L2 caches

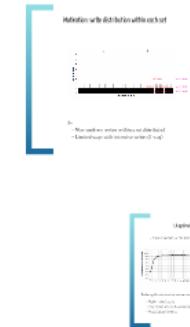
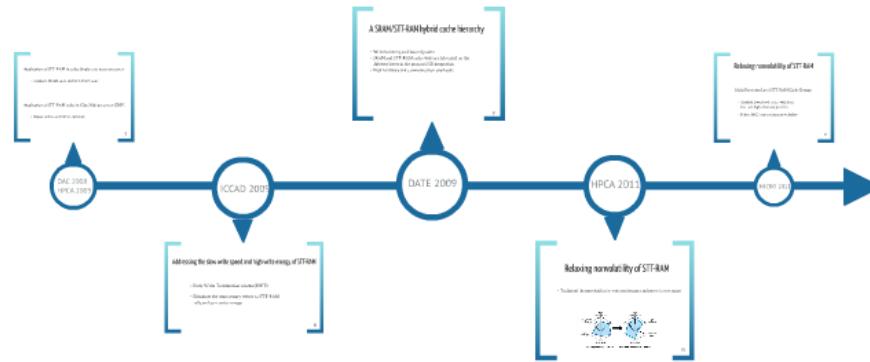
L1 caches in CMPs use write-through mechanism

- 3 level cache hierarchy



Related Works

STT-RAMs in CMP caches



Related Works

STT-RAMs in CMP caches



Application of STT-RAM in cache Single-core microprocessor

- Compare SRAM cache and STT-RAM cache

Application of STT-RAM cache in Chip Multiprocessor (CMP)

- Impact of the costly write operation

7



DAC 2008
HPCA 2009

Application of STT-RAM in cache Single-core microprocessor

- Compare SRAM cache and STT-RAM cache

Application of STT-RAM cache in Chip Multiprocessor (CMP)

- Impact of the costly write operation



ICCAD 2009

Addressing the slow write speed and high write energy of STT-RAM

- Early Write Termination scheme(EWT)
- Eliminate the unnecessary writes to STT-RAM cells and save write energy

Addressing the slow write speed and high write energy of STT-RAM

- Early Write Termination scheme(EWT)
- Eliminate the unnecessary writes to STT-RAM cells and save write energy

A SRAM/STT-RAM hybrid cache hierarchy

- Write buffering and data migration
- SRAM and STT-RAM cache ways are fabricated on the different layers in the proposed 3D integration
- High hardware and communication overheads

9

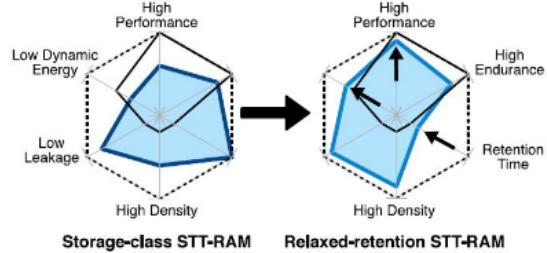


A SRAM/STT-RAM hybrid cache hierarchy

- Write buffering and data migration
- SRAM and STT-RAM cache ways are fabricated on the different layers in the proposed 3D integration
- High hardware and communication overheads

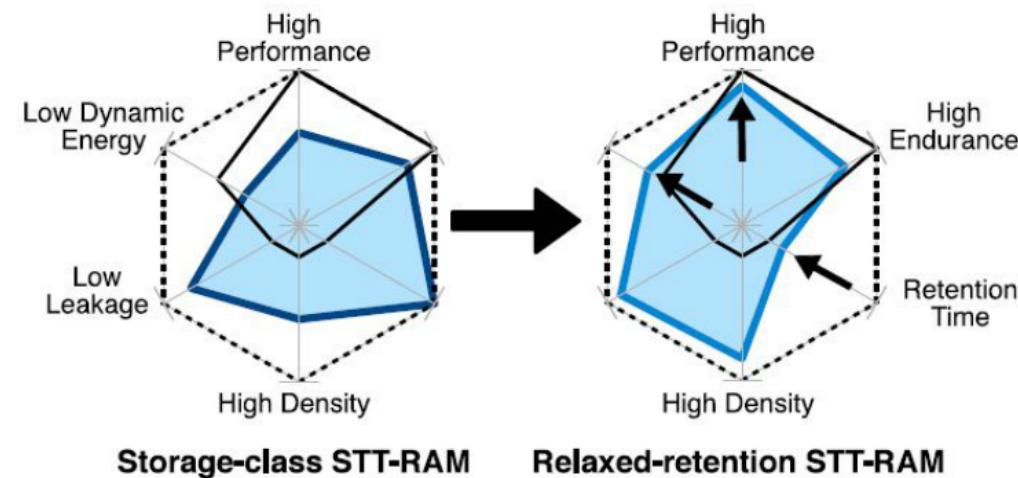
Relaxing nonvolatility of STT-RAM

- Trading off the nonvolatility for write performance and power improvement



Relaxing nonvolatility of STT-RAM

- Trading off the nonvolatility for write performance and power improvement



Relaxing nonvolatility of STT-RAM

Multi Retention Level STT-RAM Cache Designs

- A hybrid lower level cache with both low- and high-retention portions
- Reduce MTJ area to relax non-volatility

11



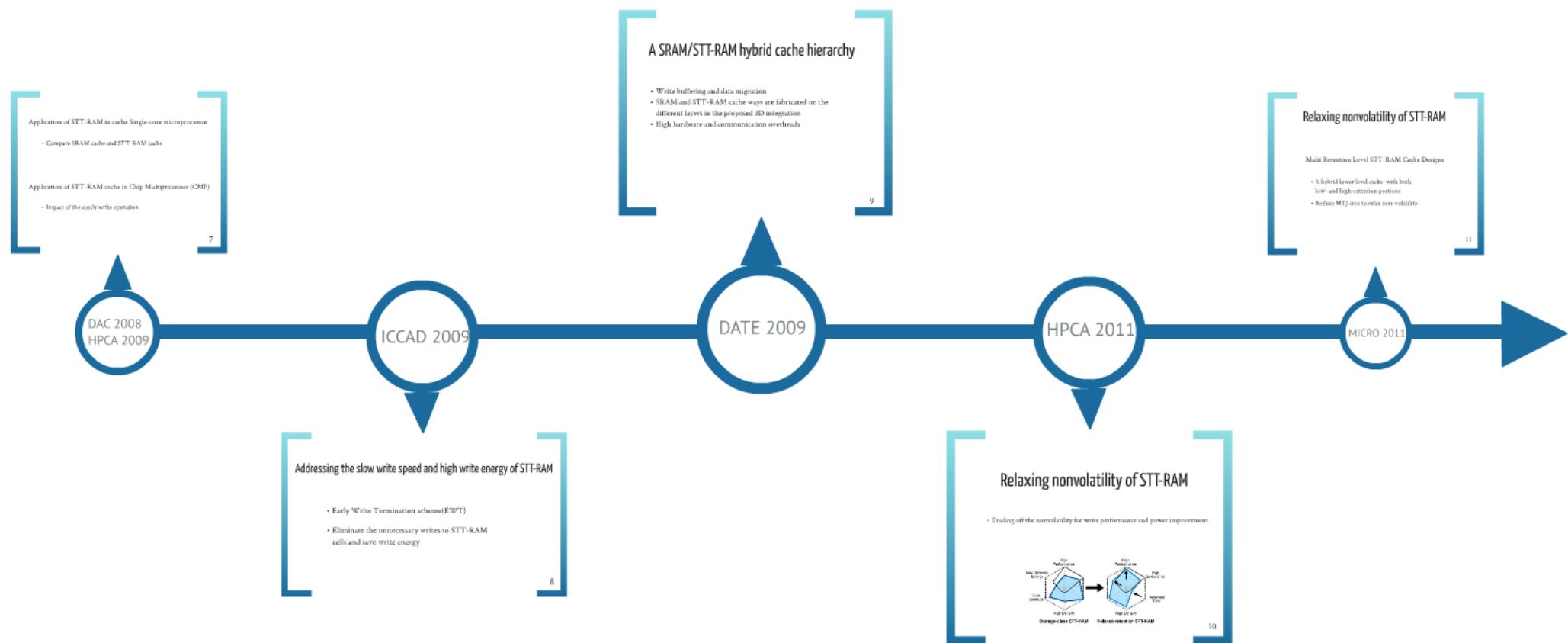
MICRO 2011

Relaxing nonvolatility of STT-RAM

Multi Retention Level STT-RAM Cache Designs

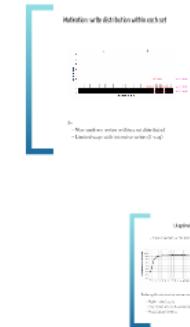
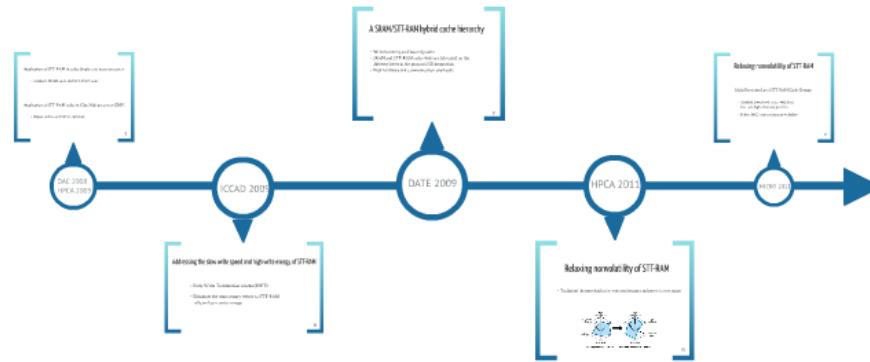
- A hybrid lower level cache with both low- and high-retention portions
- Reduce MTJ area to relax non-volatility

STT-RAMs in CMP caches



Related Works

STT-RAMs in CMP caches



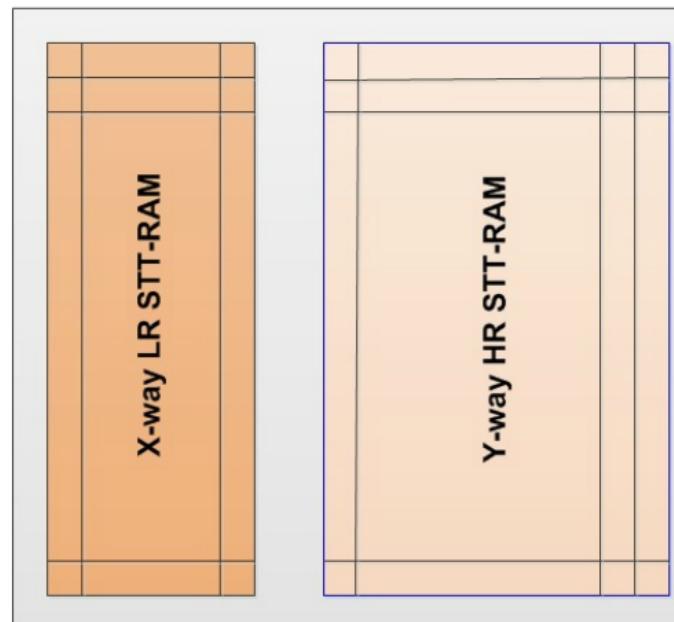
Proposed Architecture

Ideal cache requirements:

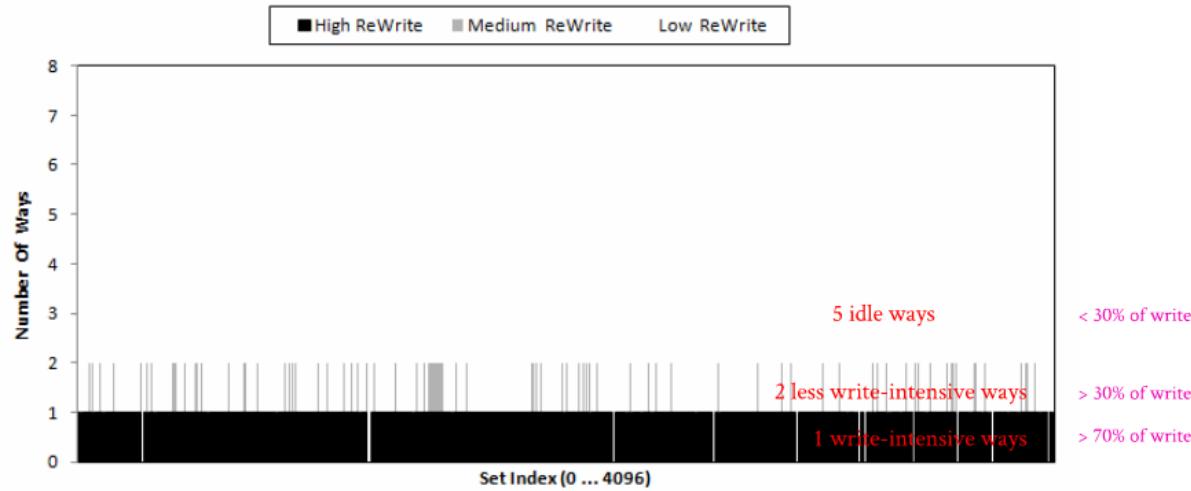
- High retention STT-RAM array
 - low leakage energy
 - non-volatility
 - low read time and energy

Ideal cache requirements:

- High retention STT-RAM array
 - low leakage energy
 - non-volatility
 - low read time and energy
- Low retention STT-RAM array
 - Low write time and energy

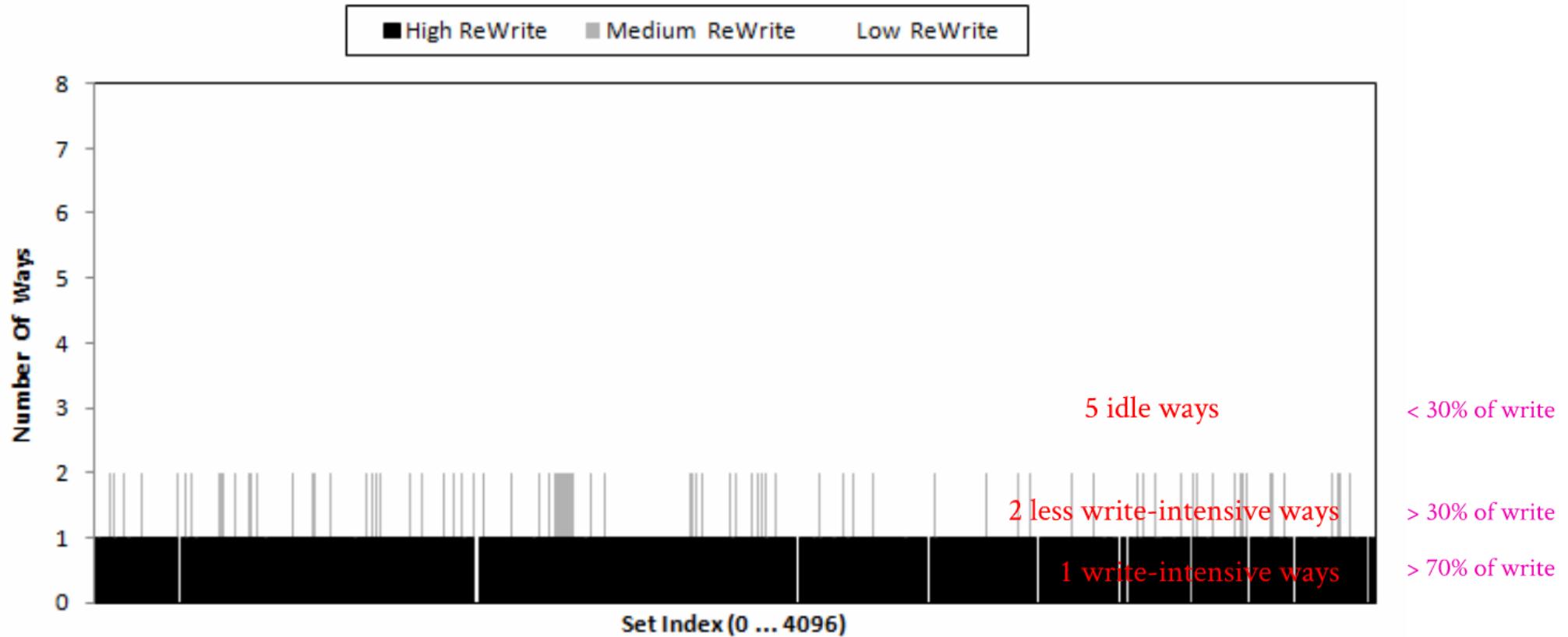


Motivation: write distribution within each set

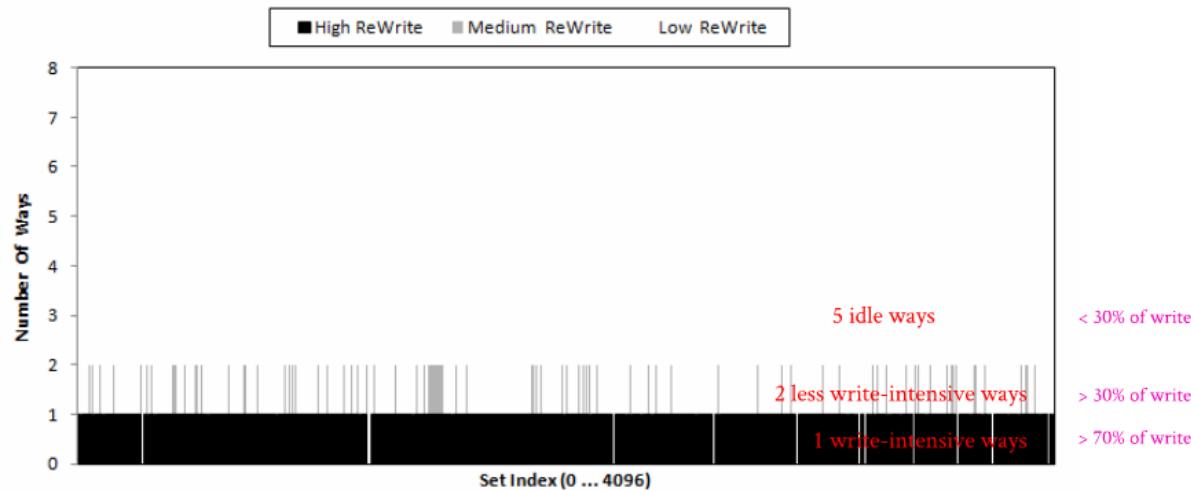


So

- Non-uniform writes within a set distributed
- Limited ways with intensive writes (1 way)



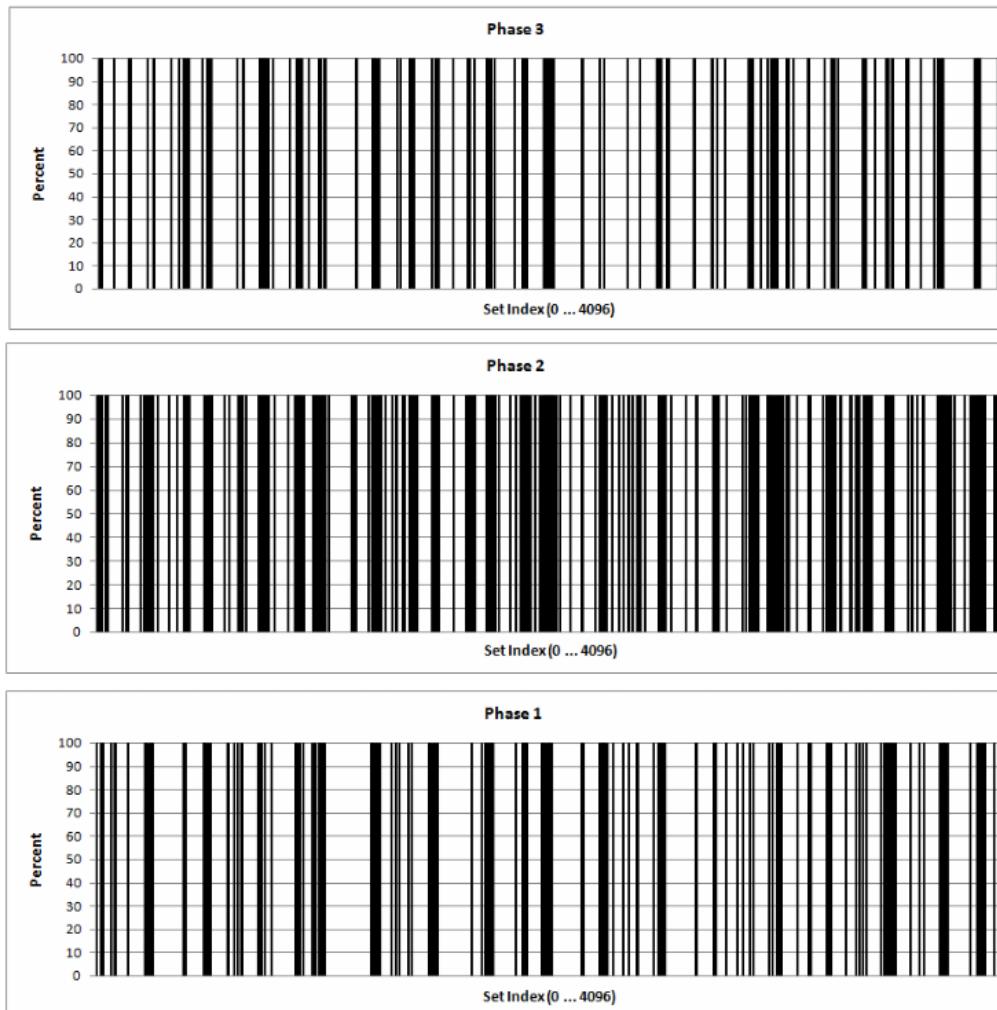
Motivation: write distribution within each set



So

- Non-uniform writes within a set distributed
- Limited ways with intensive writes (1 way)

Motivation: write distribution between different sets



Managing writes onto cache arrays

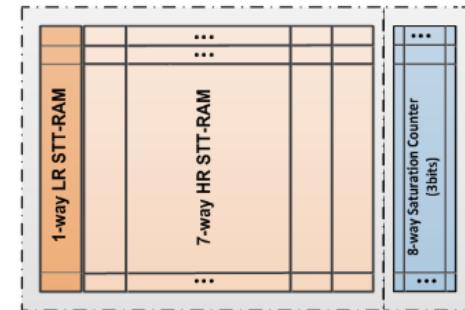
- maintain write-stressed data into the LR STT-RAM array within each set

Improve utilization of LR-region

- Intra-Set Write Management

- Per way small SCs

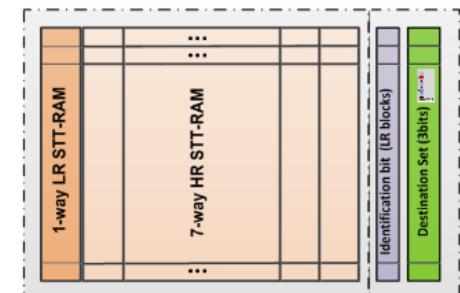
Intra-Set W rite Management



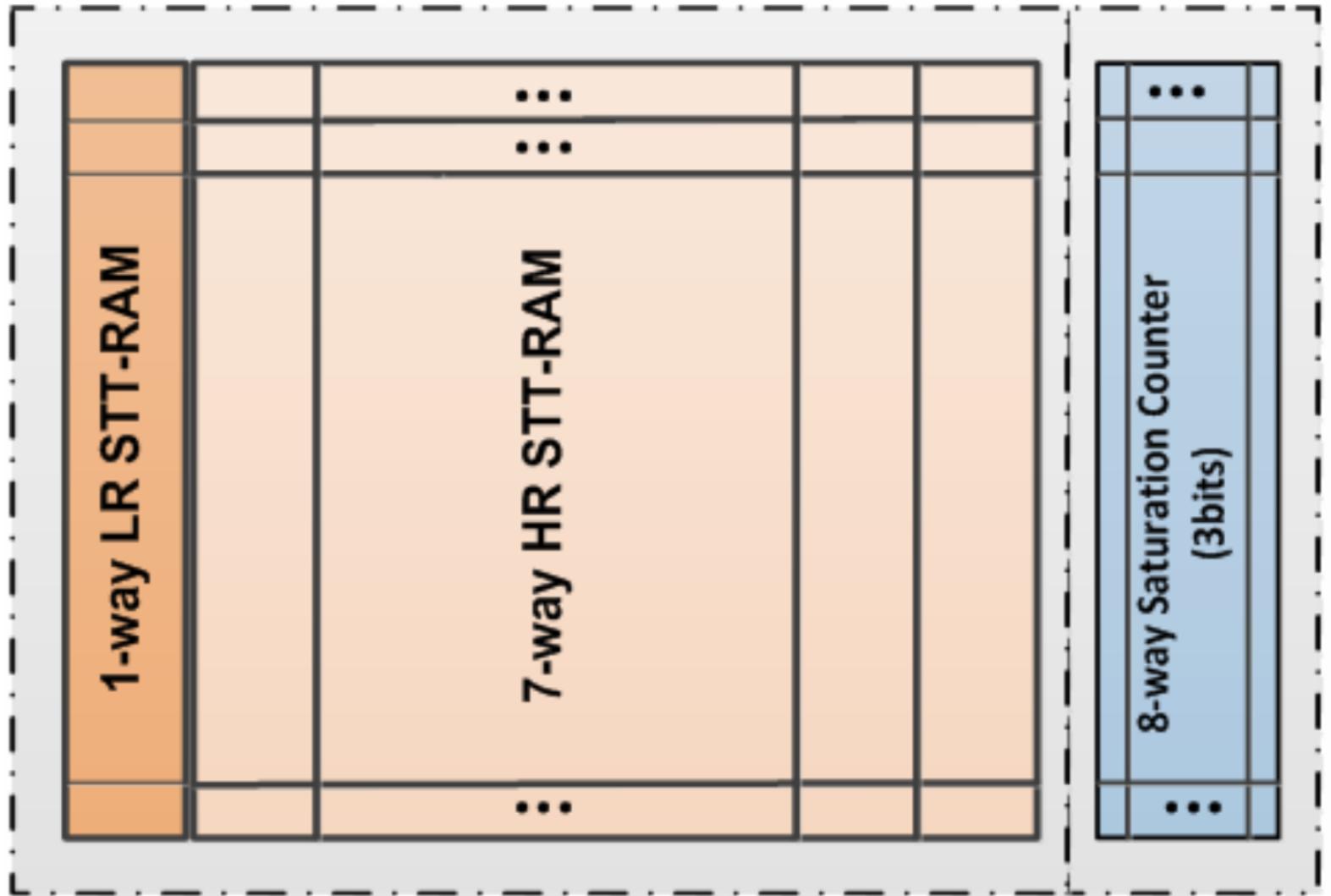
- Inter-Set Write Management

- Semi-static destination set selection
 - Keep trace of a established redirection(3bit per set)
 - Searching cache
 - Breaking a merge

Inter-Set W rite Management



Intra-Set Write Management



Managing writes onto cache arrays

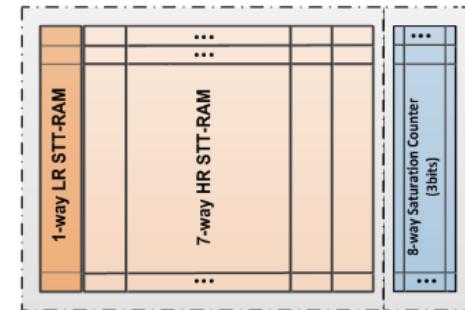
- maintain write-stressed data into the LR STT-RAM array within each set

Improve utilization of LR-region

- Intra-Set Write Management

- Per way small SCs

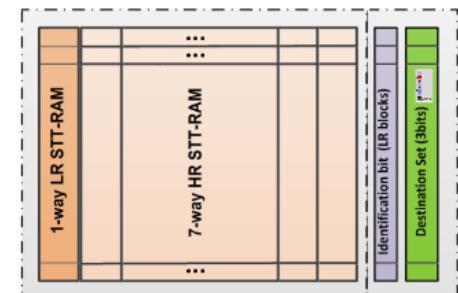
Intra-Set W rite Management



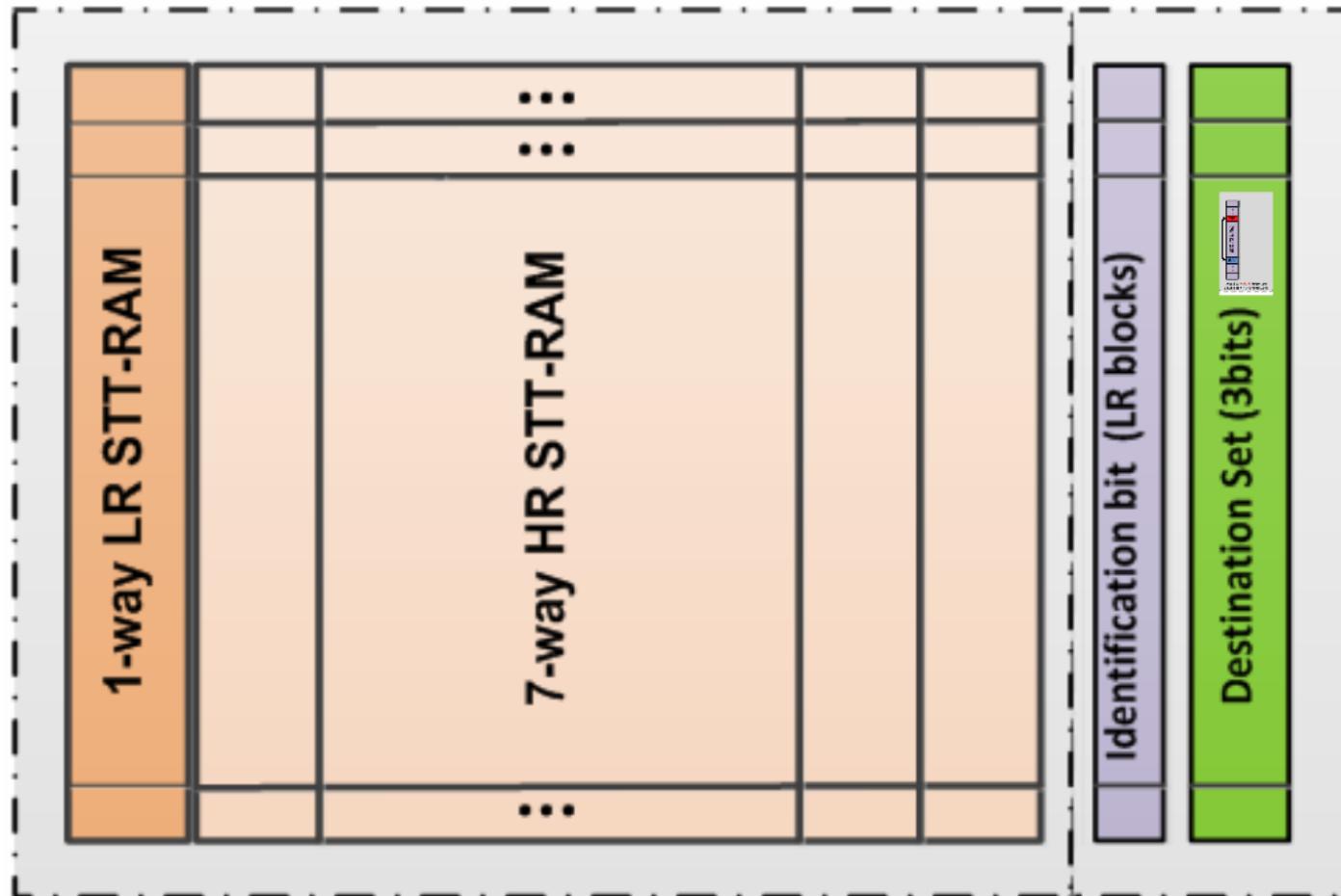
- Inter-Set Write Management

- Semi-static destination set selection
 - Keep trace of a established redirection(3bit per set)
 - Searching cache
 - Breaking a merge

Inter-Set W rite Management



Inter-Set Write Management

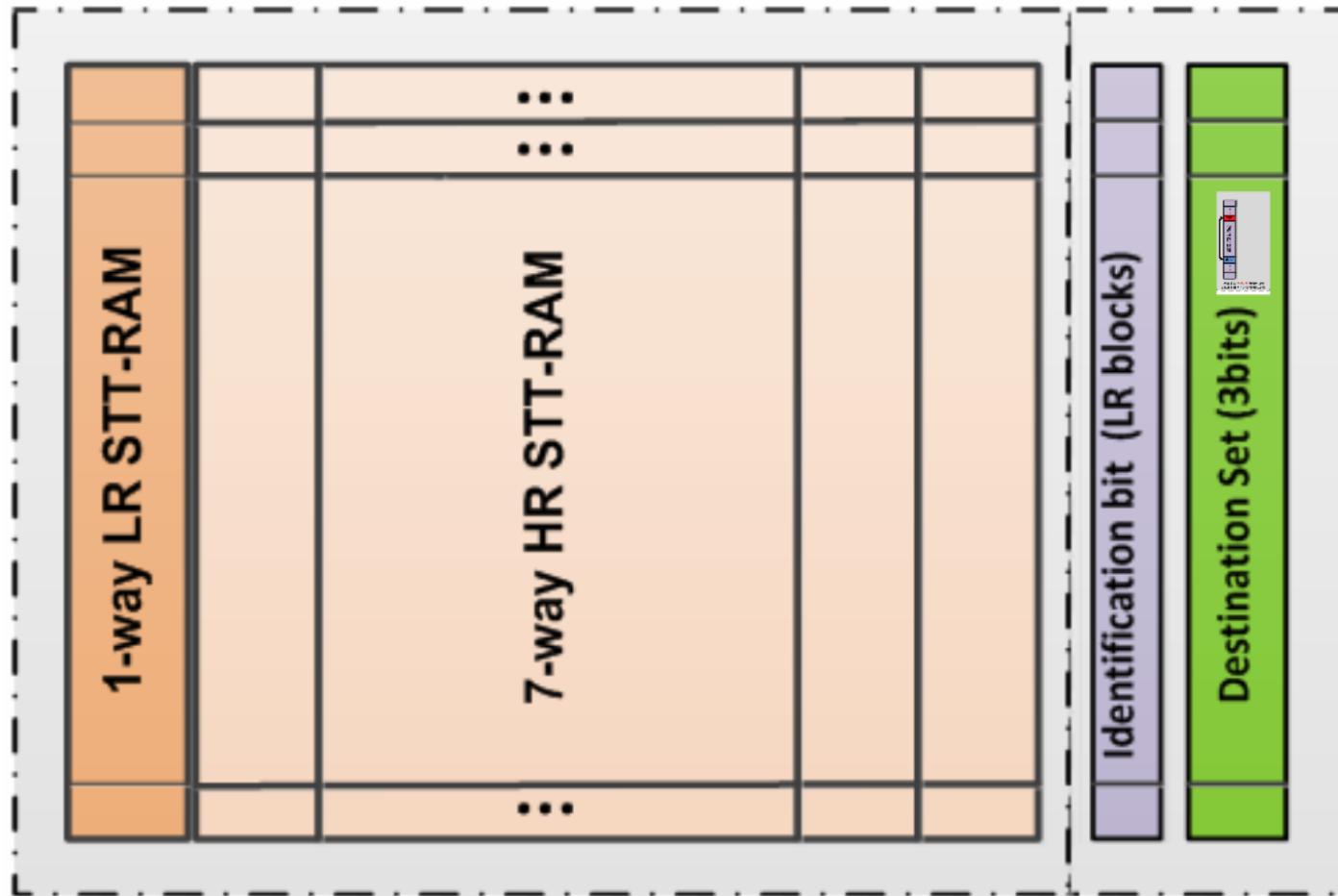




Hot Set = **011**11010...011

Cold Set = **101**11010...011

Inter-Set Write Management



Managing writes onto cache arrays

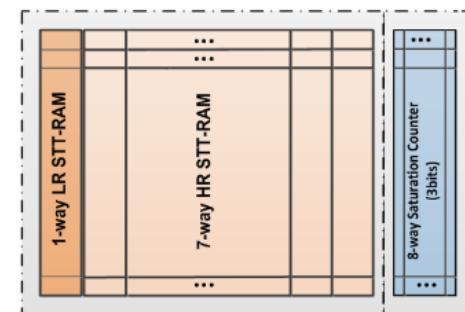
- maintain write-stressed data into the LR STT-RAM array within each set

Improve utilization of LR-region

- Intra-Set Write Management

- Per way small SCs

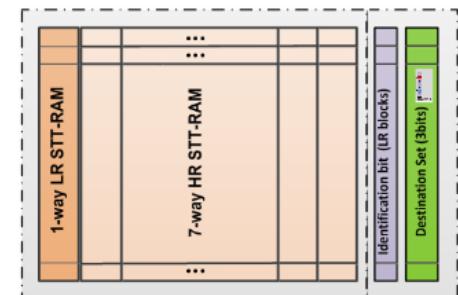
Intra-Set W rite Management



- Inter-Set Write Management

- Semi-static destination set selection
 - Keep trace of a established redirection(3bit per set)
 - Searching cache
 - Breaking a merge

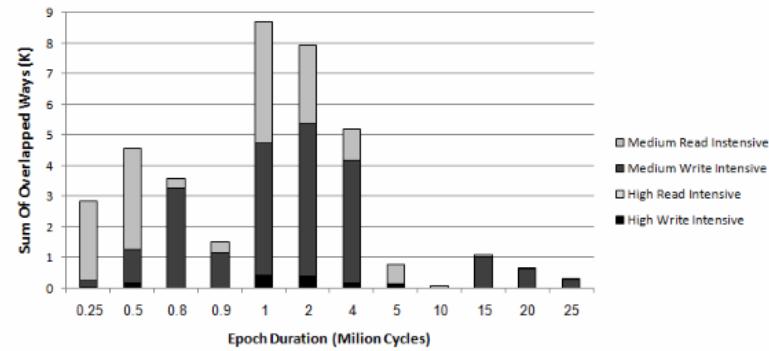
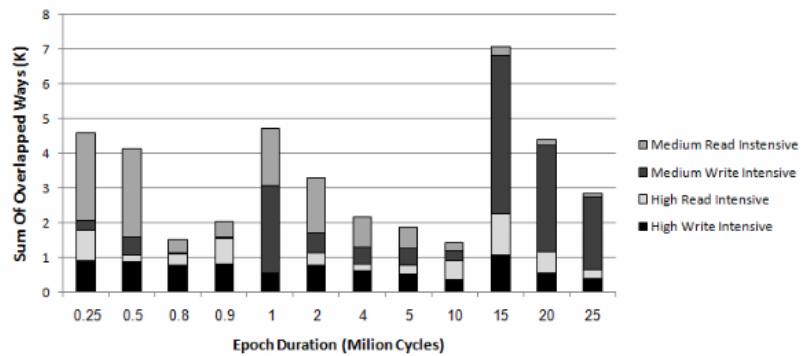
Inter-Set W rite Management

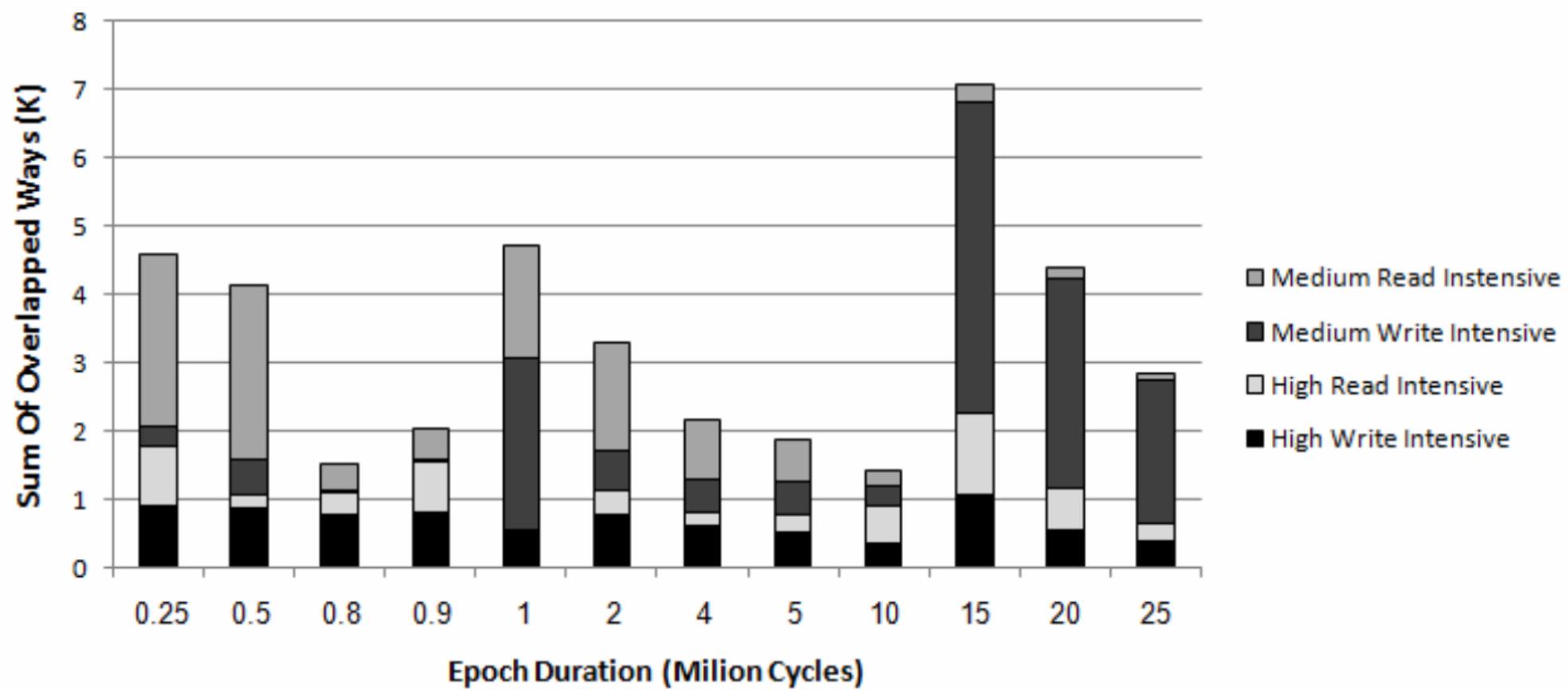


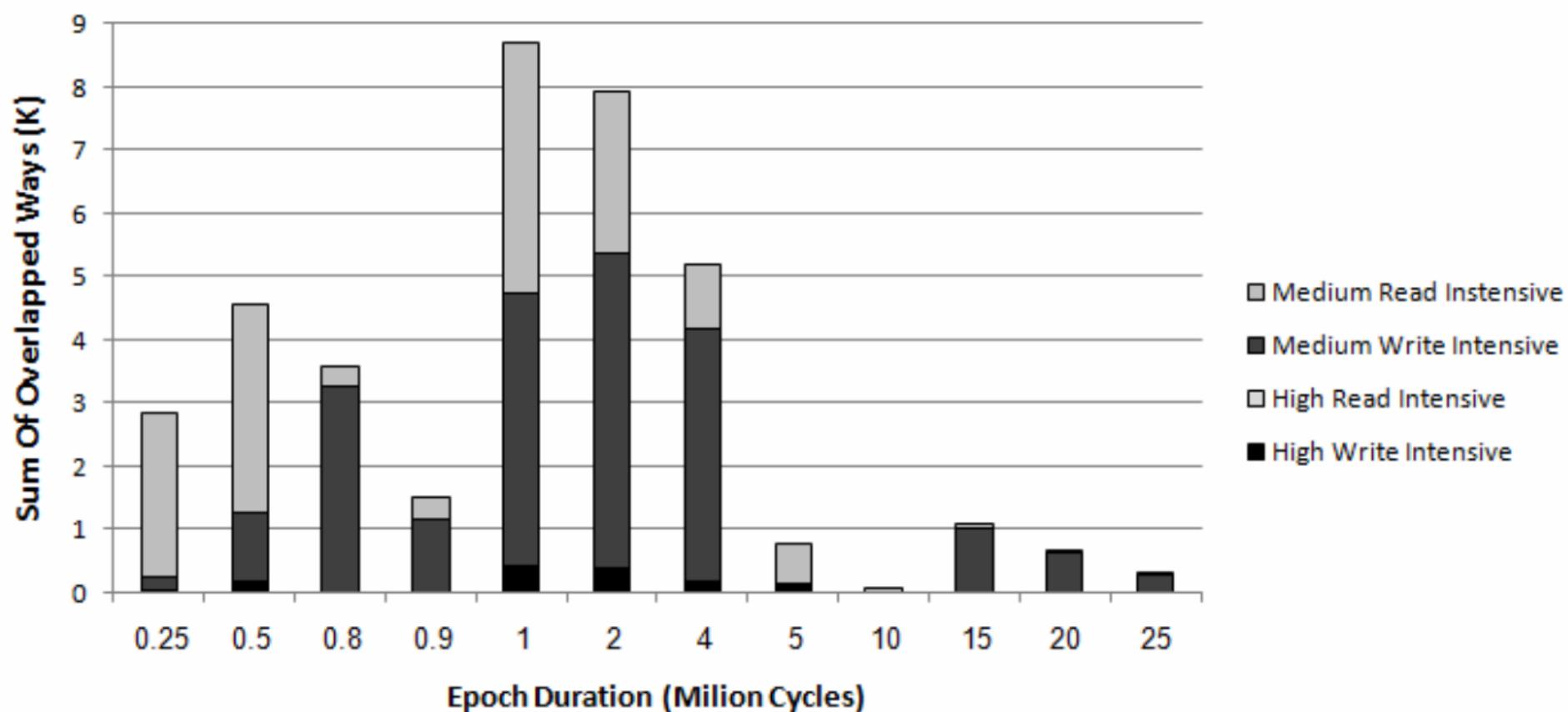
Best time interval

At the beginning of optimum time interval :

- Redirect W-intensive blocks from HR-region to LR-region in a set



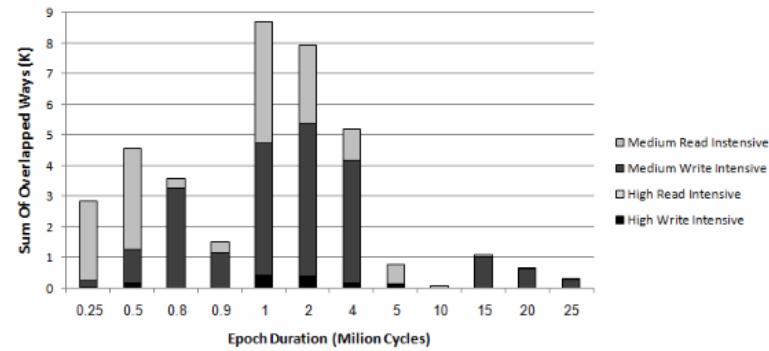
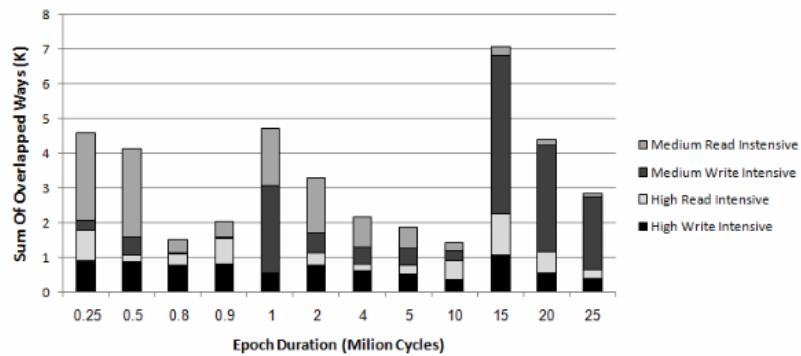




Best time interval

At the beginning of optimum time interval :

- Redirect W-intensive blocks from HR-region to LR-region in a set



latency and energy of LR & HR cache



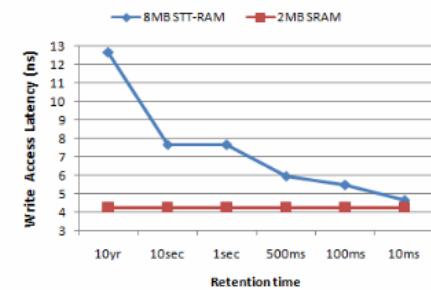
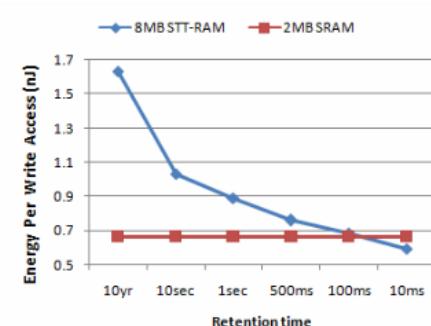
$$T_{store} = 1\text{ns} \cdot e^{\Delta}$$

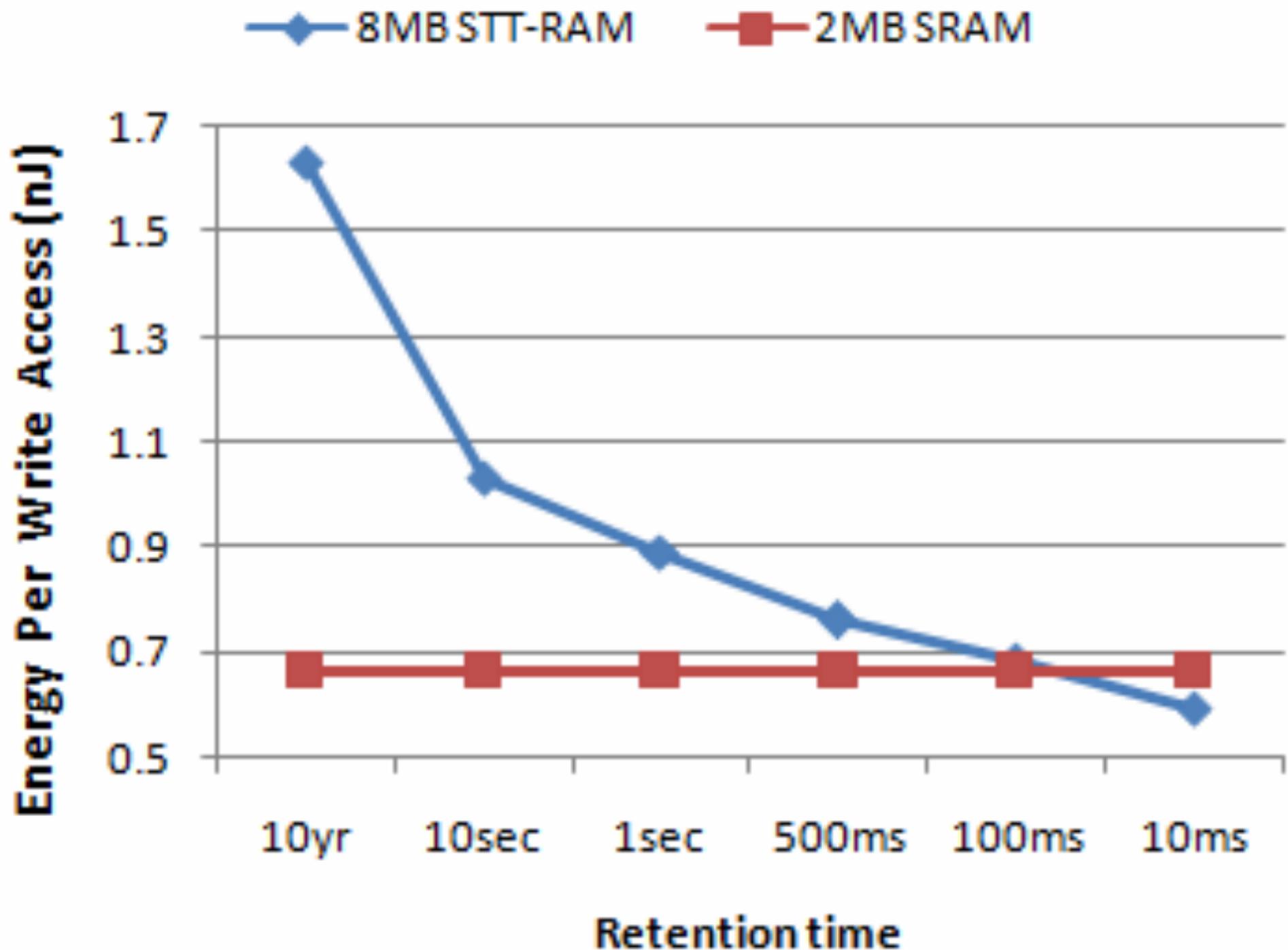
$$\Delta = \frac{M_s H_k A t_F}{2k_B T}$$

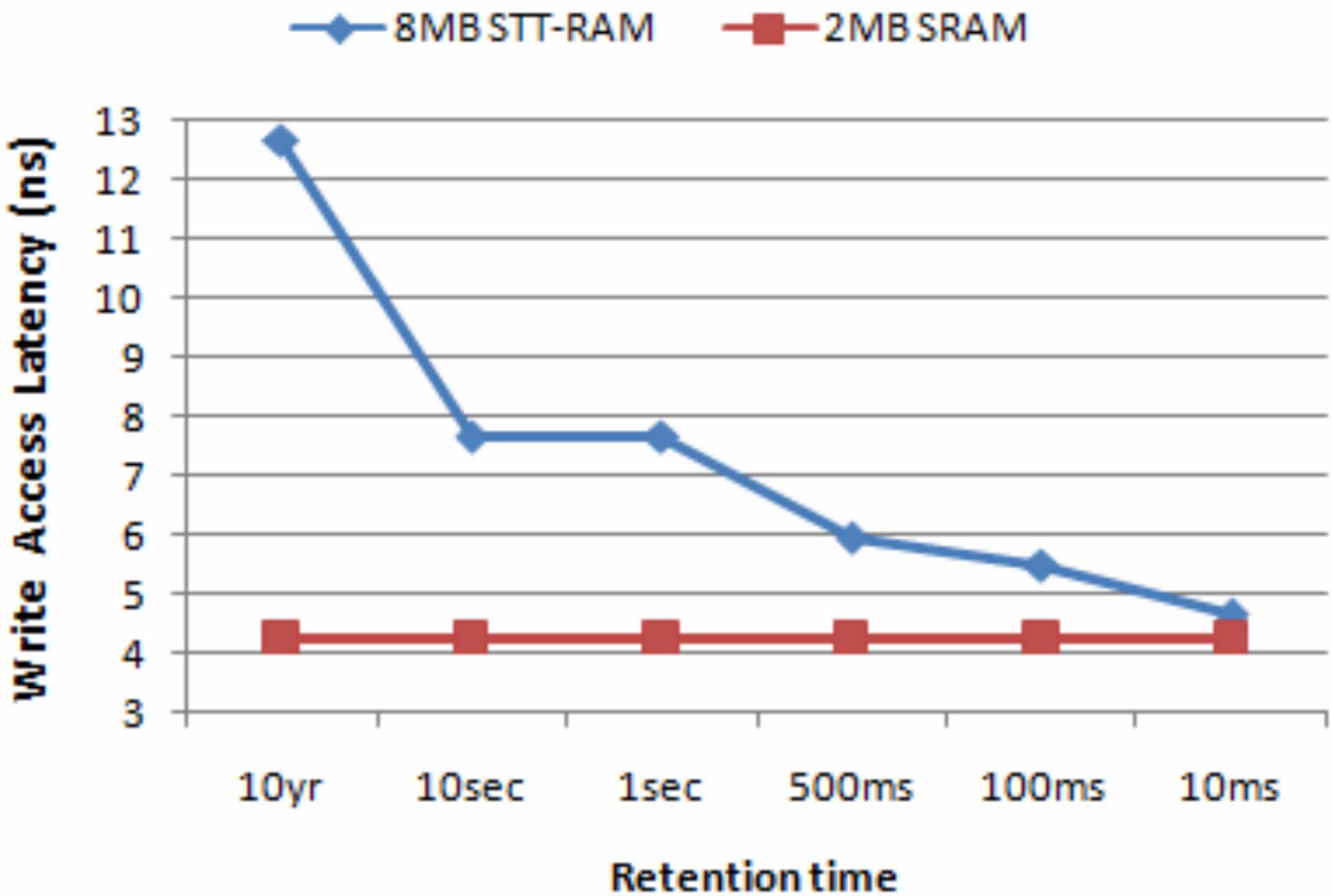
$$J_{c0} = \left(\frac{2e}{h}\right) \left(\frac{\alpha}{\eta}\right) (t_F M_s) (H_k \pm H_{ext} + 2\pi M_s)$$

$$I_c(\tau) = A \times J_c(\tau)$$

$$Write\ Energy = V_{Write} \times I_c(\tau) \times \tau$$







latency and energy of LR & HR cache



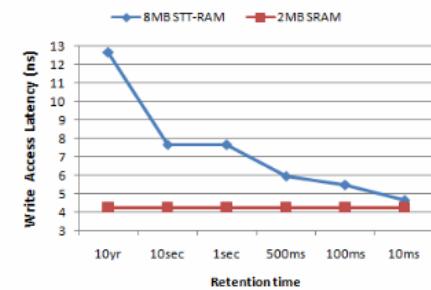
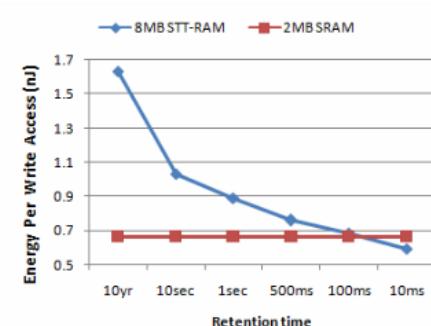
$$T_{store} = 1\text{ns} \cdot e^{\Delta}$$

$$\Delta = \frac{M_s H_k A t_F}{2k_B T}$$

$$J_{c0} = \left(\frac{2e}{h}\right) \left(\frac{\alpha}{\eta}\right) (t_F M_s) (H_k \pm H_{ext} + 2\pi M_s)$$

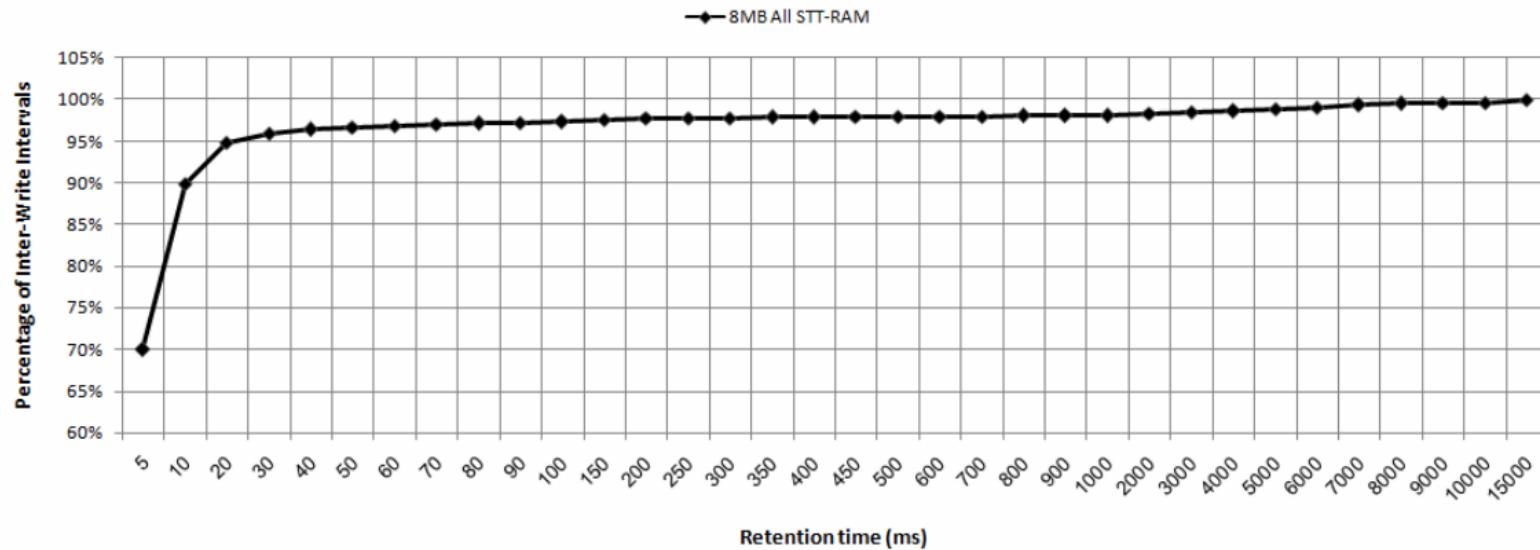
$$I_c(\tau) = A \times J_c(\tau)$$

$$Write\ Energy = V_{Write} \times I_c(\tau) \times \tau$$



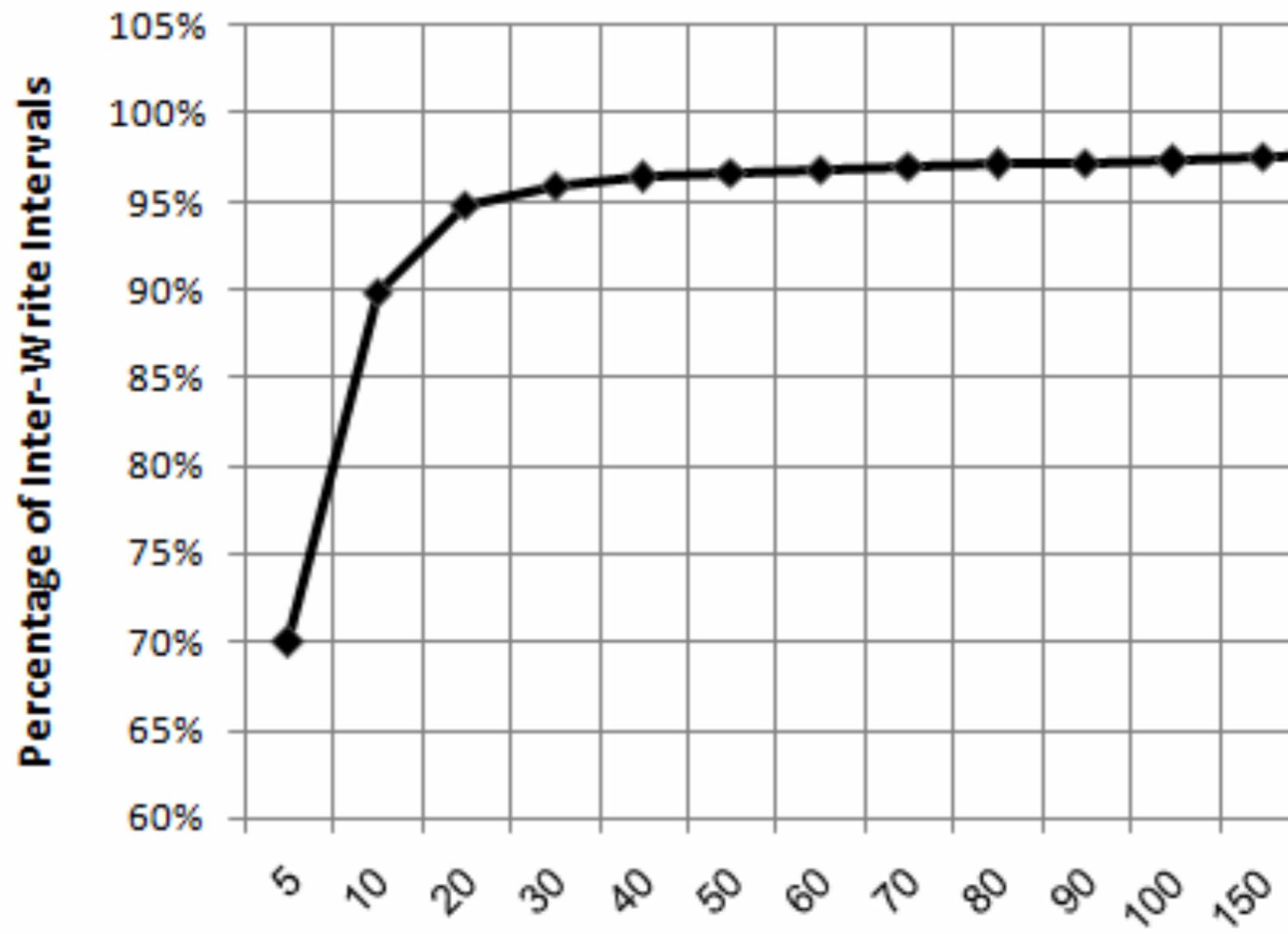
LR optimal retention time

- More than 90% of the inter-write time intervals are smaller than 10ms



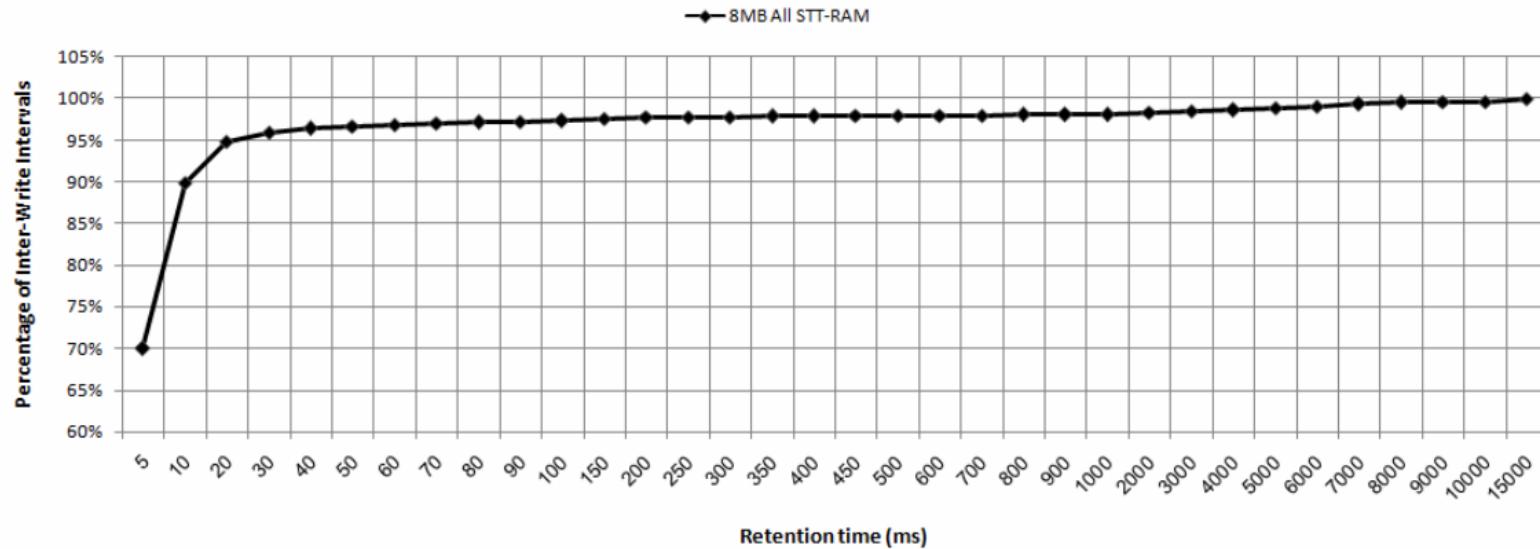
Reducing the retention time to micro second scale:

- Highly volatile cache
- Any refresh scheme becomes impractical for the large lower level cache.
- Degraded performance



LR optimal retention time

- More than 90% of the inter-write time intervals are smaller than 10ms

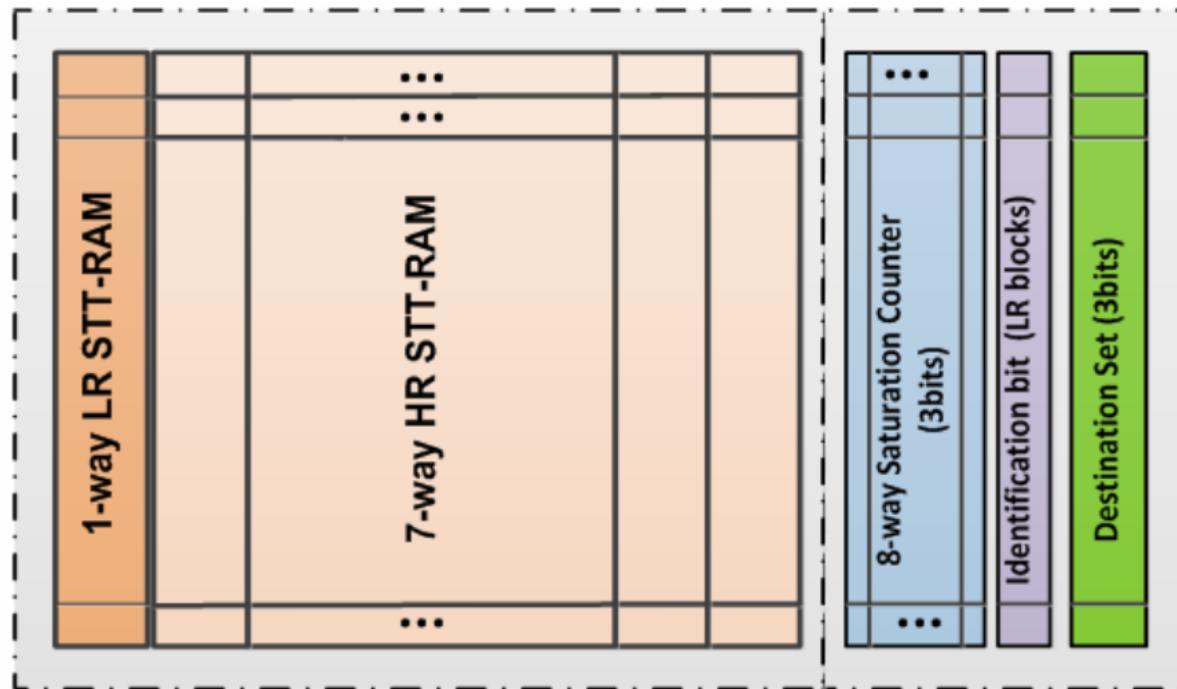


Reducing the retention time to micro second scale:

- Highly volatile cache
- Any refresh scheme becomes impractical for the large lower level cache.
- Degraded performance

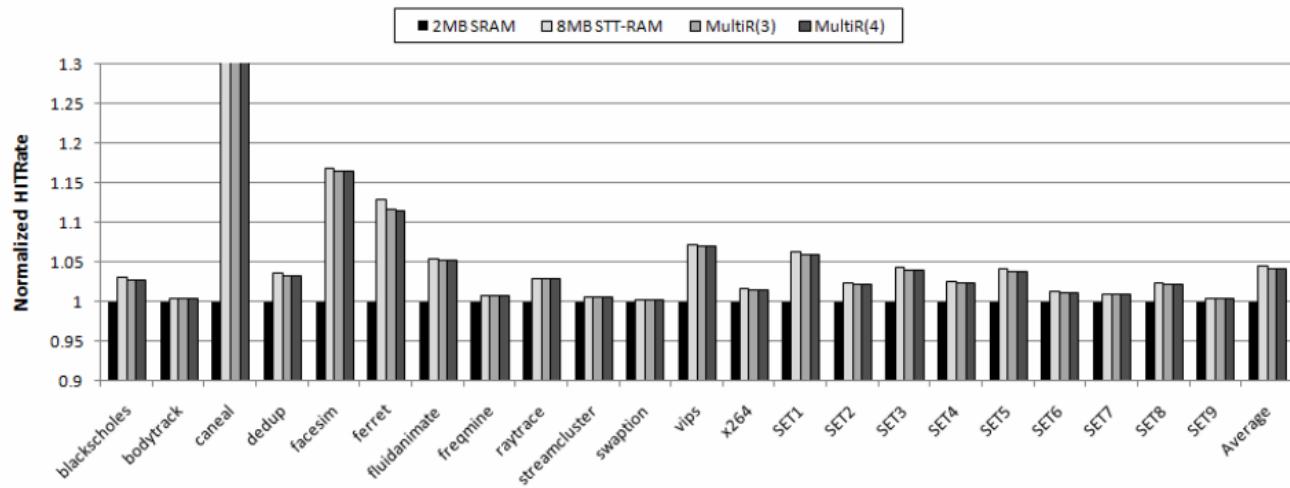
Redundancy for each set :

- One Saturation Counter for each way in HR-region (w-intensive block)
- 1 bit for each way in LR-region (valid bit)
- 3 bits destination set identification for each set

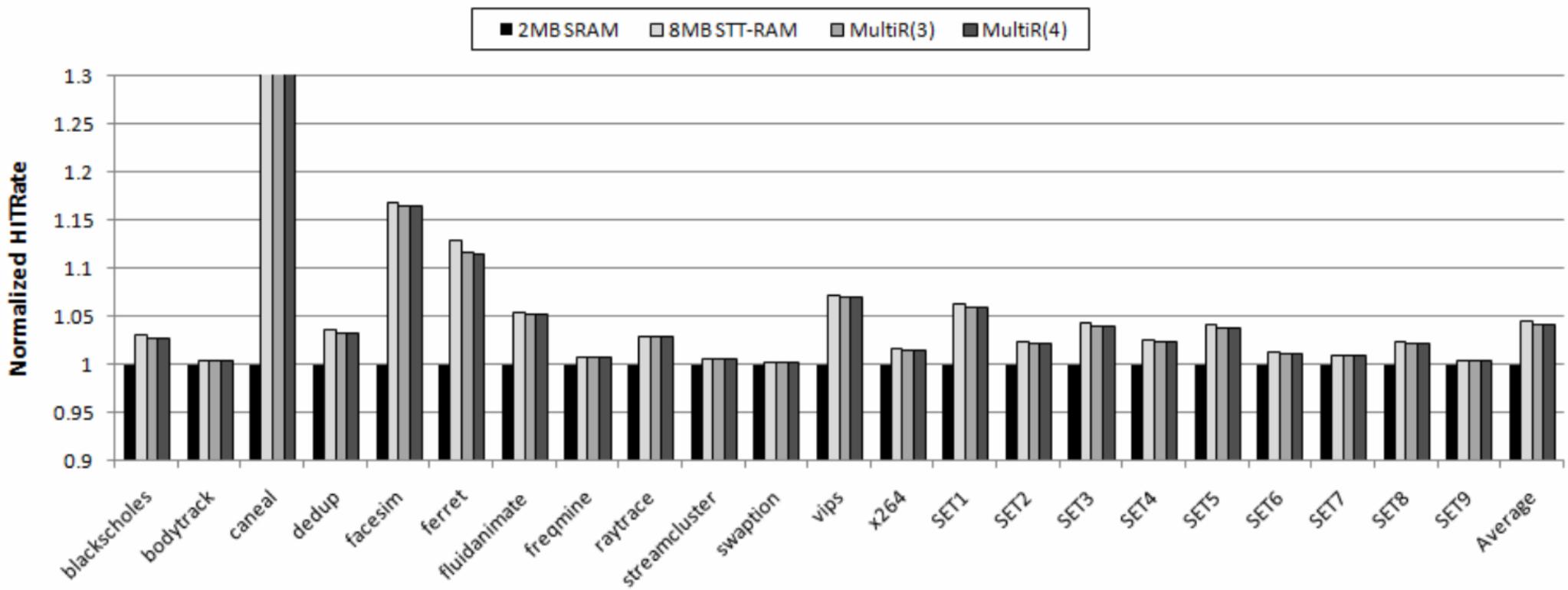


Refreshing mechanism for data blocks in LR region

- Maintain a copy of LR block in one HR block in each set
- Refresh overhead :

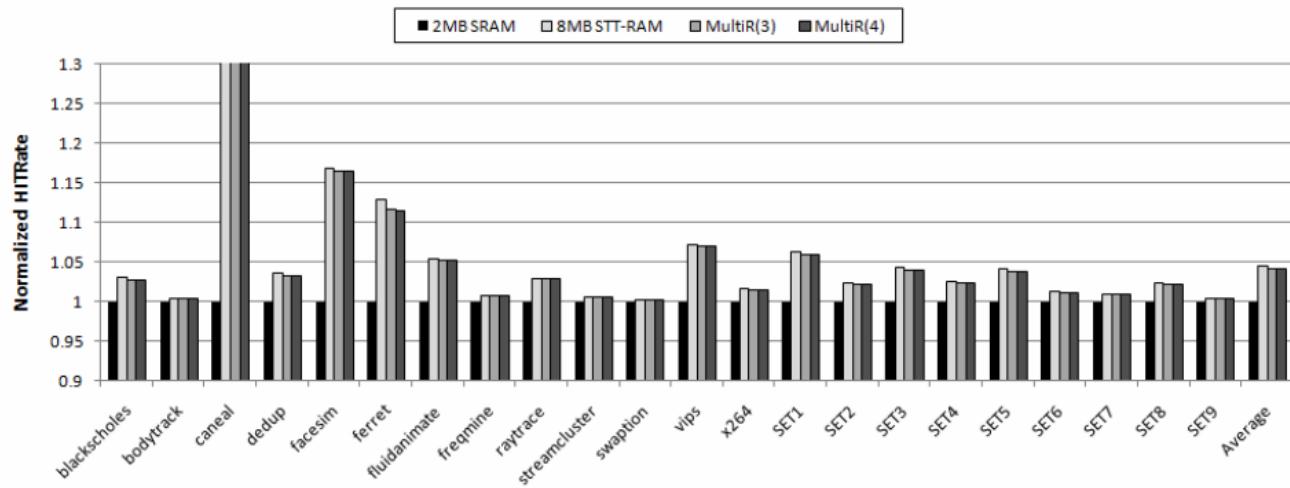


- Maintain a copy of LR block in one HR block in each set
- Refresh overhead :



Refreshing mechanism for data blocks in LR region

- Maintain a copy of LR block in one HR block in each set
- Refresh overhead :



Proposed Architecture



Ideal cache requirements:

- High utilization (75% to 80% usage)
- Low miss rate
- Low read time and energy

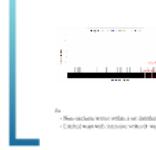
Low utilization (75% to 80% usage)

- Low access time and energy



13

Migration with distribution within each set



14

Migration with distribution between different sets



15

Best time interval



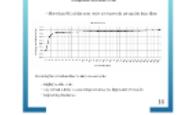
16

Always empty of old blocks



17

Updated writer file



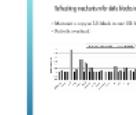
18

Migration for S set



19

Replacing each byte for data blocks in S register



20

Evaluation

Simulation Platform

System-level simulation

- Virtutech SIMOS full-system simulator
- 4 cores
- Cycle accurate
- + CACTI timing model rounded by system latency
- + CACTI energy model
- Benchmark: ROI of the SPEC CPU2006 + PARSEC-2 programs

Simulator configuration:

Processor	Intel Xeon E5-2620 v2	Processor clock	2.1 GHz
L1 Cache	32 KB per core, 8 MB total, 4-way associativity, 128B cache line size, 10 ns access time	Latency	10 ns
L2 Cache	128 KB shared, 16-way associativity, 128B cache line size, 10 ns access time	Latency	10 ns
L3 Cache	16 MB shared, 16-way associativity, 128B cache line size, 10 ns access time	Latency	10 ns
Memory	4 GB DDR3-1600, 16-bit wide memory bus	Latency	10 ns

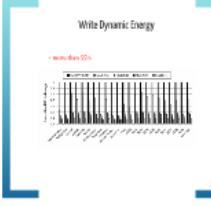
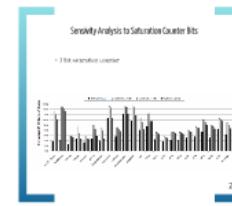
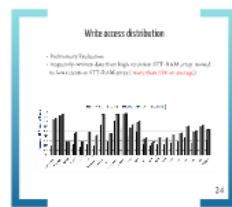
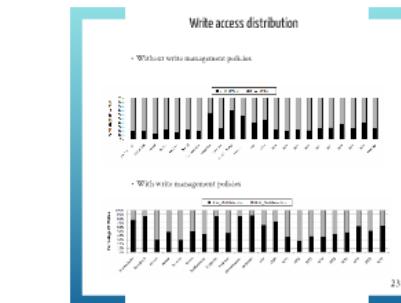
21

Design approaches

+ 1.3 2MB SRAM replacement with 8MB STT-RAM with same area

Cache Configuration	Op.	Execution Time	Dramatic Energy	Power	Area
SRAM 1.3M	R	—	0.0001 J/cycle	0.0001 W	1.3 MB
SRAM 8MB	R	—	0.0001 J/cycle	0.0001 W	8 MB
SRAM MultiR	LR, R	—	0.0001 J/cycle	0.0001 W	8 MB
SRAM MultiR	LR, R	10 ns	0.0001 J/cycle	0.0001 W	8 MB
SRAM MultiR	LR, R	10 ns	0.0001 J/cycle	0.0001 W	8 MB
SRAM MultiR	LR, R	10 ns	0.0001 J/cycle	0.0001 W	8 MB
SRAM MultiR	LR, R	10 ns	0.0001 J/cycle	0.0001 W	8 MB
SRAM MultiR	LR, R	10 ns	0.0001 J/cycle	0.0001 W	8 MB

22



Evaluation

Design approaches

Simulation Platform

System-level simulation:

- Vitutech SIMICS full-system simulator
- + Msimulator
 - Cycle accurate
- + CACTI timing model rounded by system frequency
- + CACTI energy model
- Benchmarks: ROI of the SPEC CPU2k6 + PARSEC-2 programs

Simulator configuration:

Processor	4-core SPARC-v9 core, 2.5 GHz, Running Solaris 10
L1 Cache	Split I and D cache; 4KB private; 4-way set associative; 32B line size; LRU; write-through policy; 1 port; MOESI directory-based coherency; SRAM with 1.45ns (4 cycles) access time
L2 Cache	UCA 128KB shared; 4-way set associative; 64B line size; LRU; writeback policy; SRAM bank: 2.58ns (7 cycles) latency
L3 Cache	UCA 2MB shared; 8-way set associative; 128B line size; LRU; writeback policy; SRAM bank: 4.23n (11 cycles) latency
Main Memory	4GB DRAM; 50ns (200 cycles) access latency

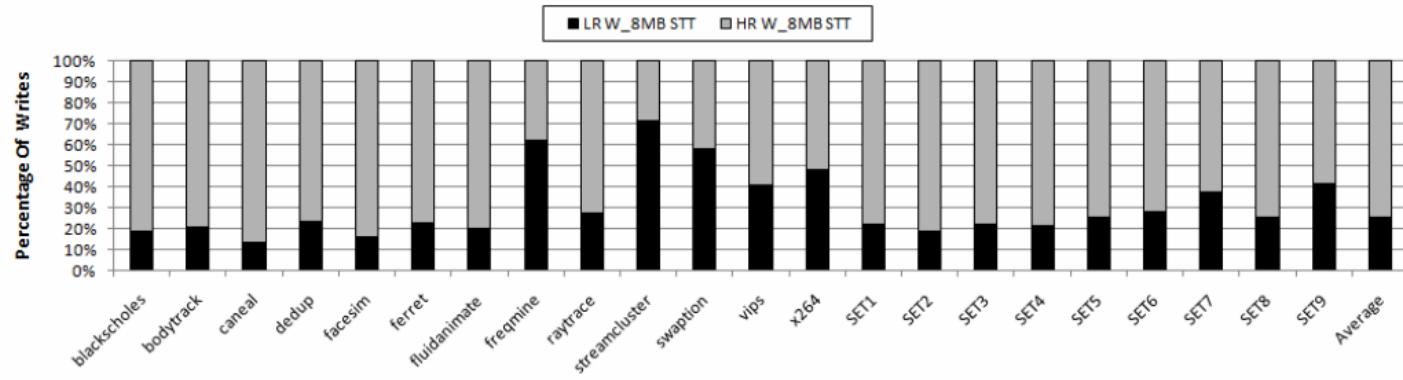
Design approaches

- L3 2MB SRAM replacement with 8MB STT-RAM with same area

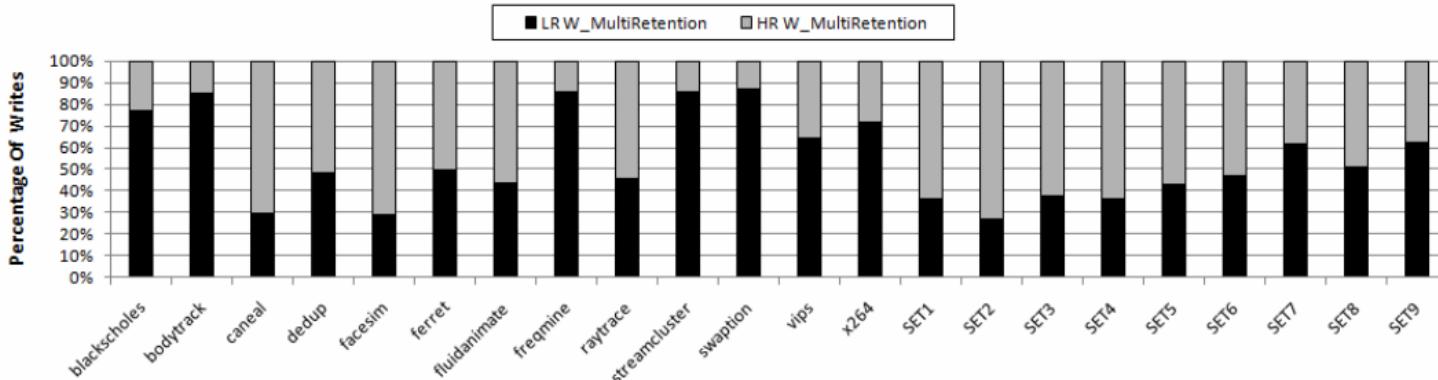
Cache Config.	Op.	Retention Time	Dynamic Energy	Pulse Duration	Total latency (+Tag)
8MB STT-RAM	R	---	0.530 nJ/access	---	4.283 ns
	W	10 year	1.635 nJ/access	10ns	12.67 ns
8MB MuIR(1)	R	---	0.530 nJ/access	---	4.283 ns
	LR-W	10ms	0.593 nJ/access	2ns	4.672 ns
	HR-W	1year	1.635 nJ/access	10ns	12.67 ns
8MB MuIR(2)	LR-W	100ms	0.682 nJ/access	2.8ns	5.472 ns
	HR-W	1sec	0.887 nJ/access	5ns	7.672 ns
8MB MuIR(3)	LR-W	10ms	0.593 nJ/access	2ns	4.672 ns
	HR-W	1sec	0.887 nJ/access	5ns	7.672 ns
8MB MuIR(4)	LR-W	10ms	0.593 nJ/access	2 ns	4.672 ns
	HR-W	100ms	0.682 nJ/access	2.8 ns	5.472 ns

Write access distribution

- Without write management policies

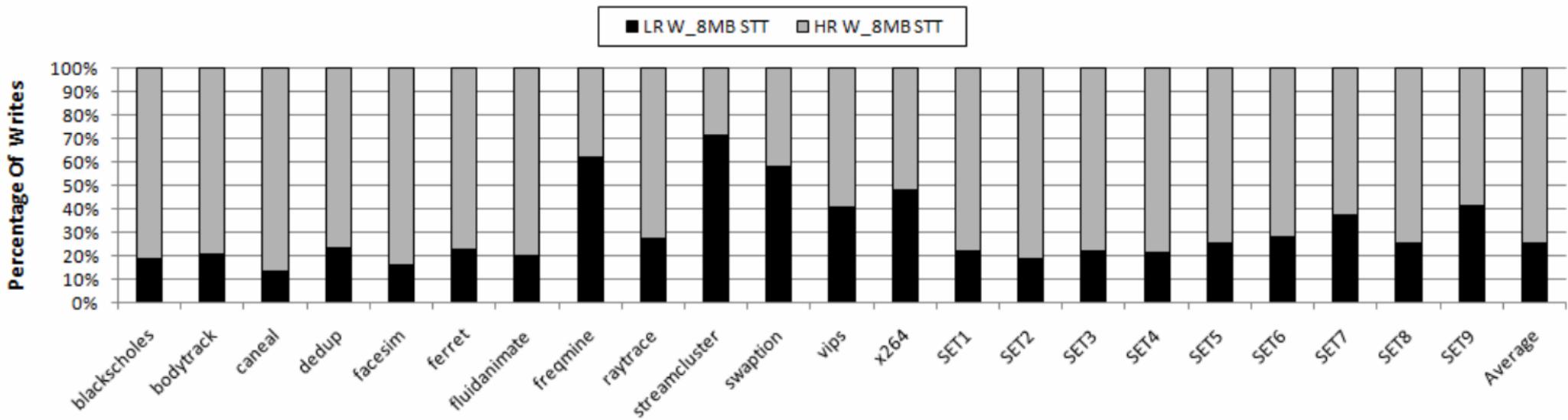


- With write management policies

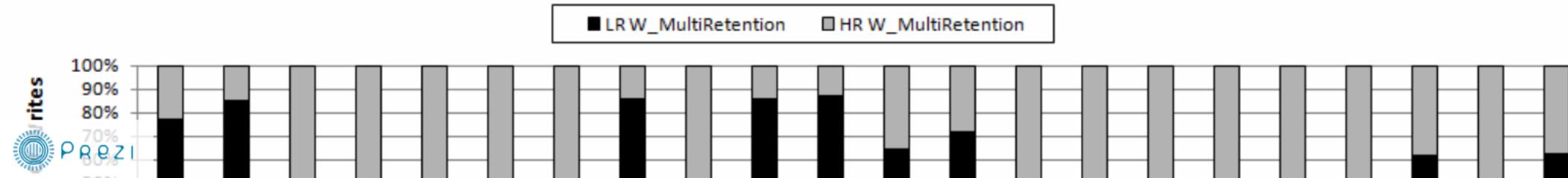


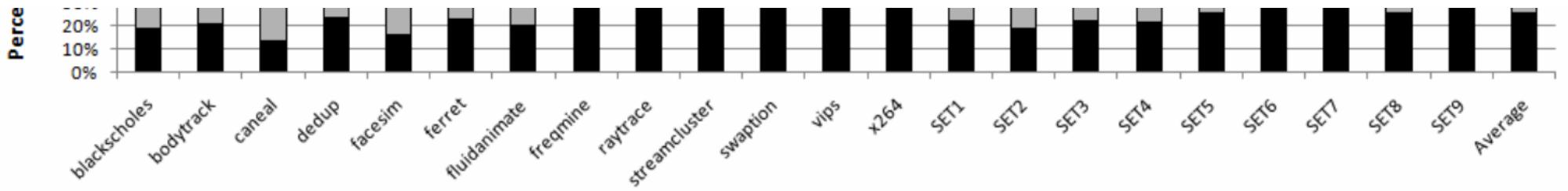
Write access distribution

- Without write management policies

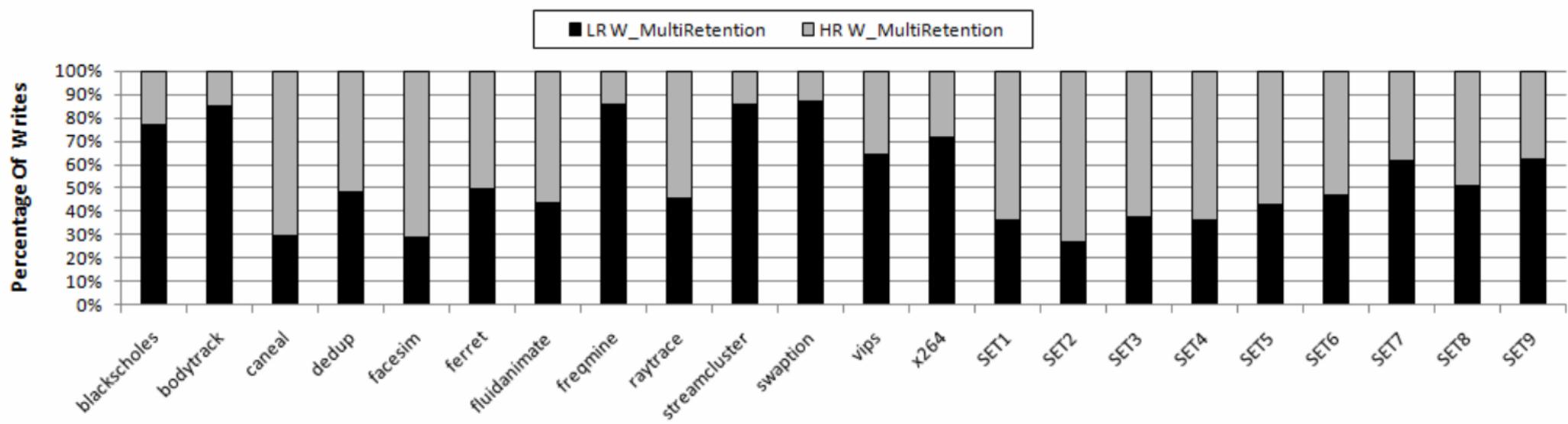


- With write management policies



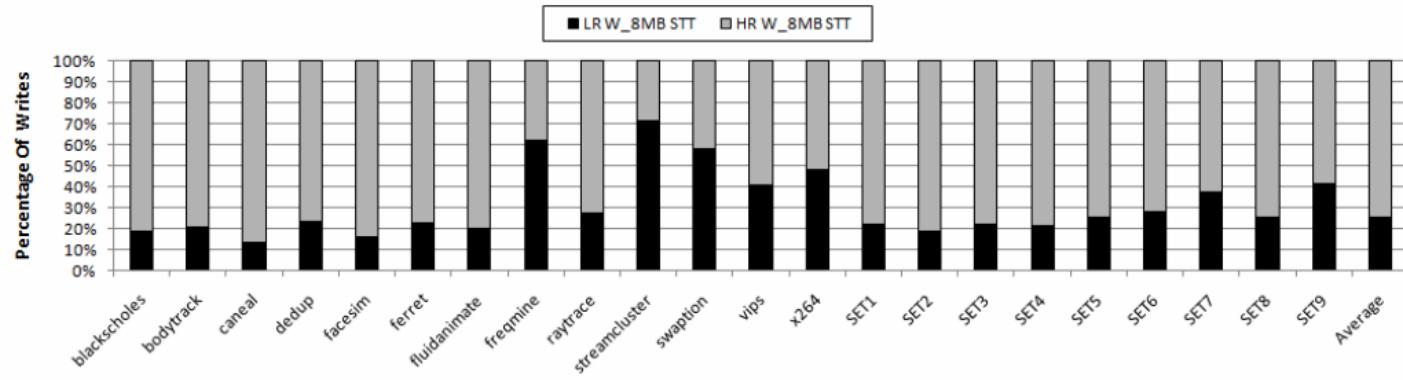


- With write management policies

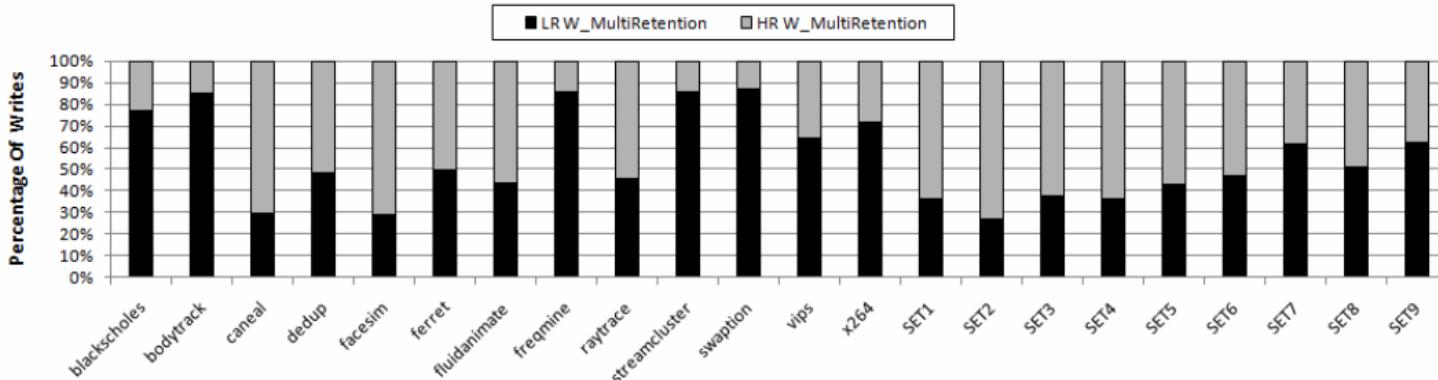


Write access distribution

- Without write management policies

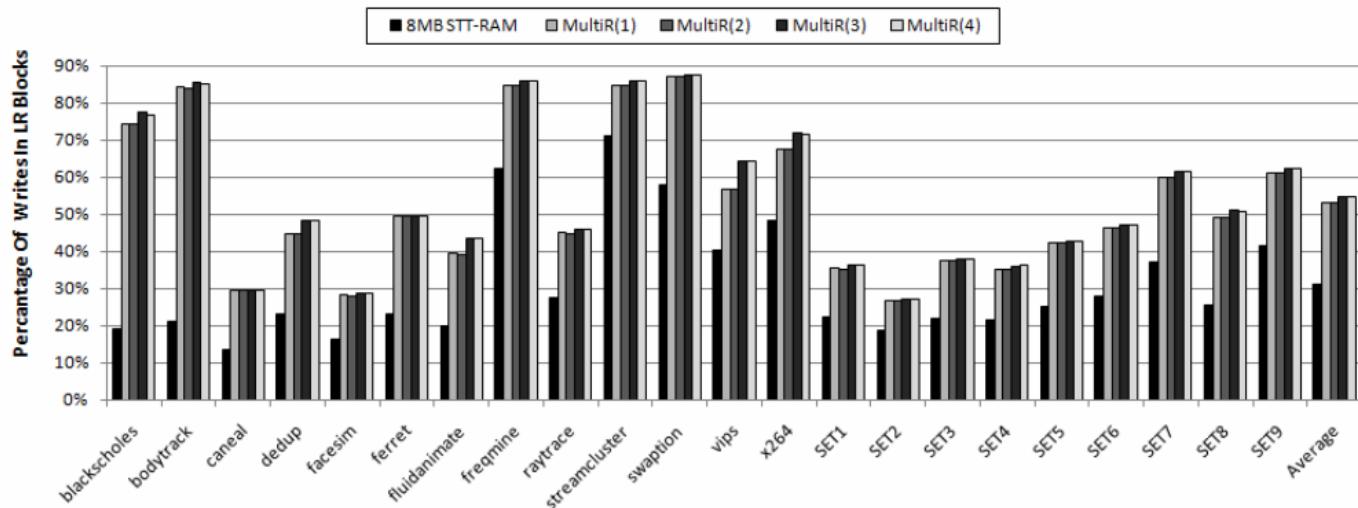


- With write management policies

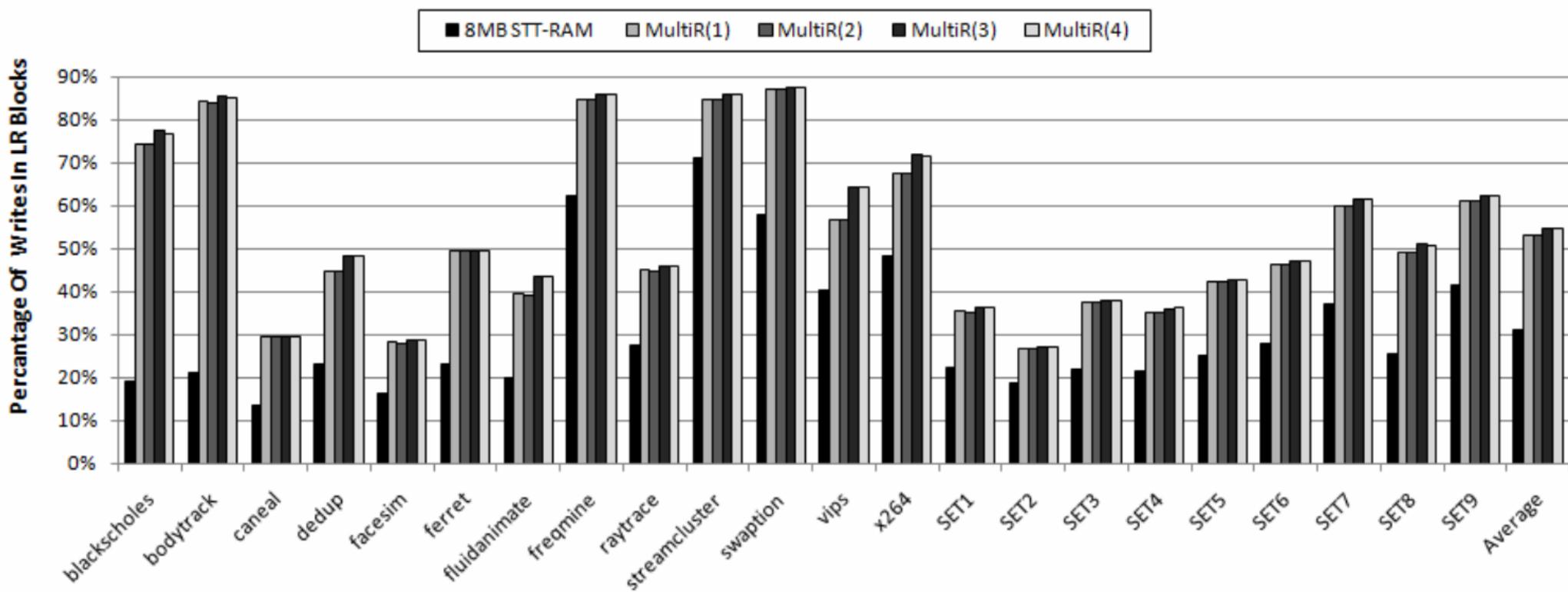


Write access distribution

- Preliminary Evaluation
- frequently-written data from high retention STT-RAM array moved to low retention STT-RAM array (**more than 55% on average**)

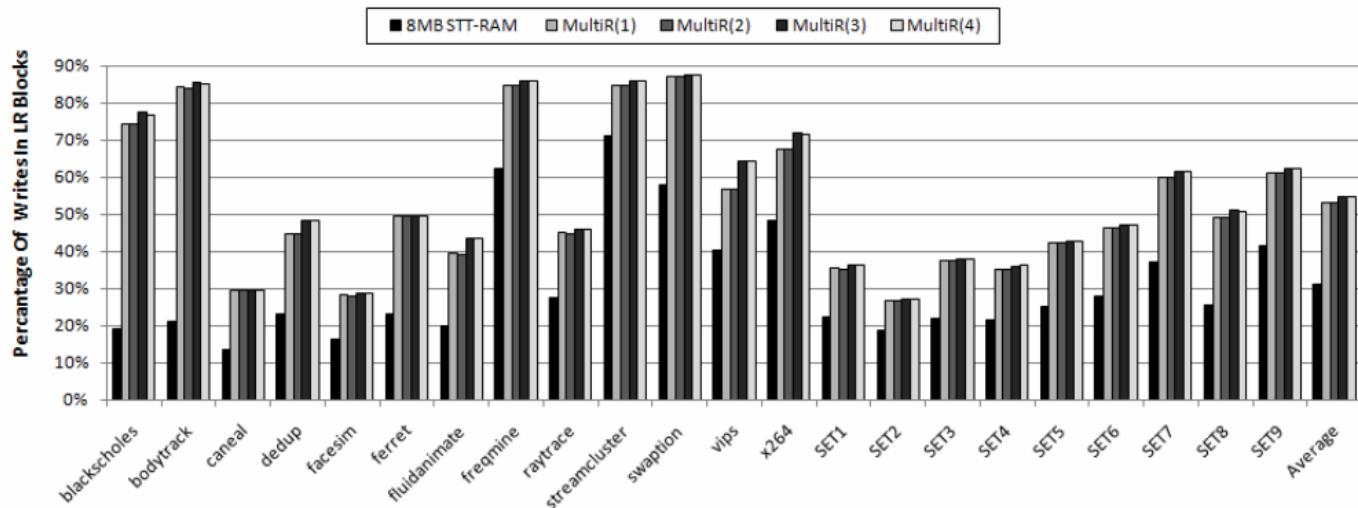


to low retention STT-RAM array (more than 55% on average)



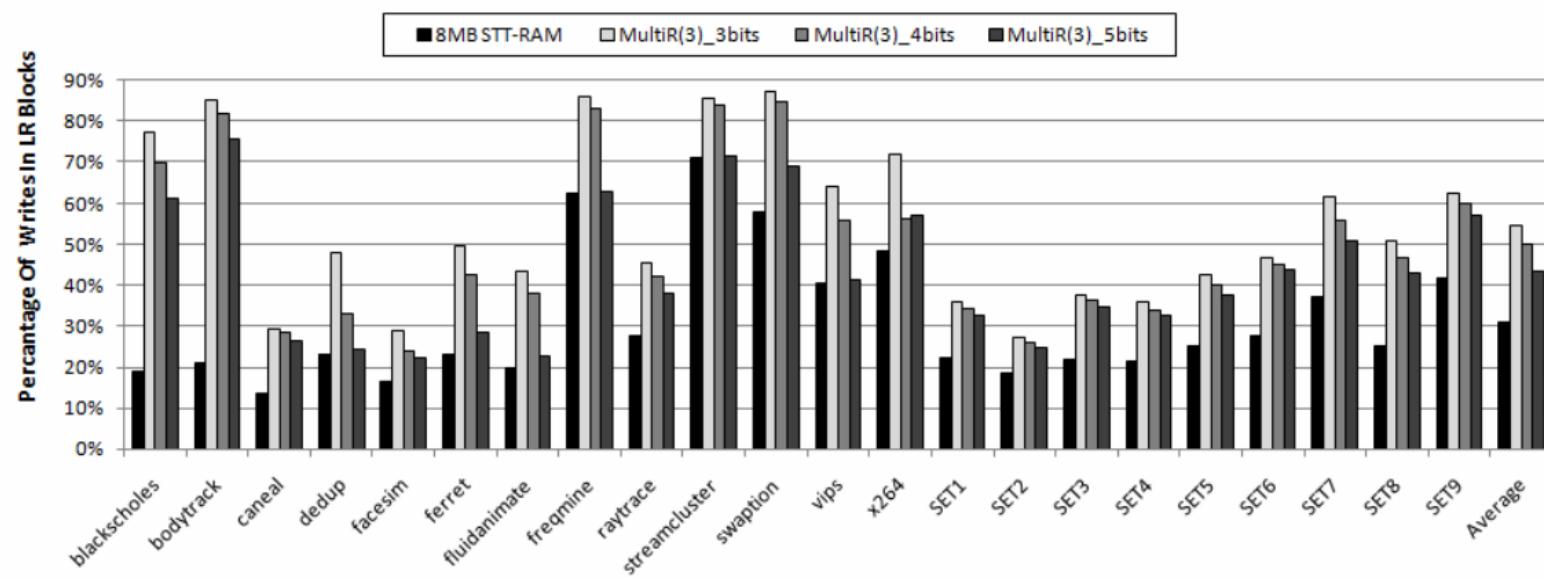
Write access distribution

- Preliminary Evaluation
- frequently-written data from high retention STT-RAM array moved to low retention STT-RAM array (**more than 55% on average**)

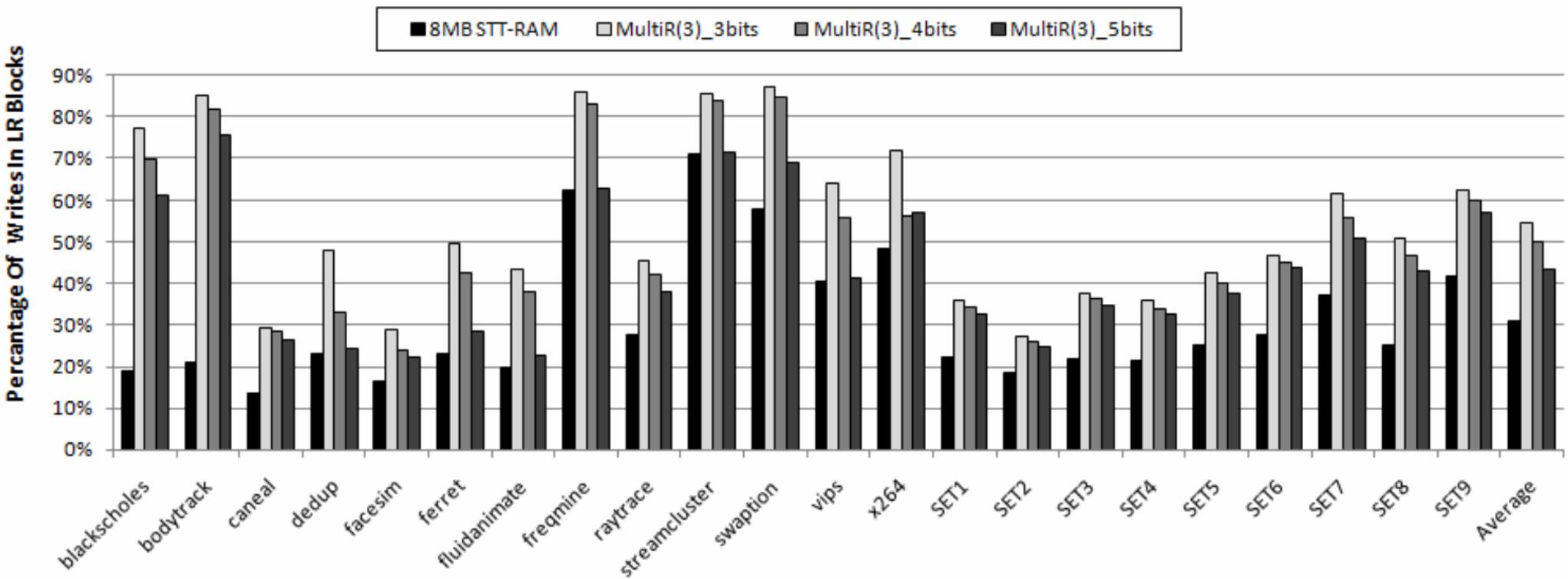


Sensitivity Analysis to Saturation Counter Bits

- 3 bit saturation counter

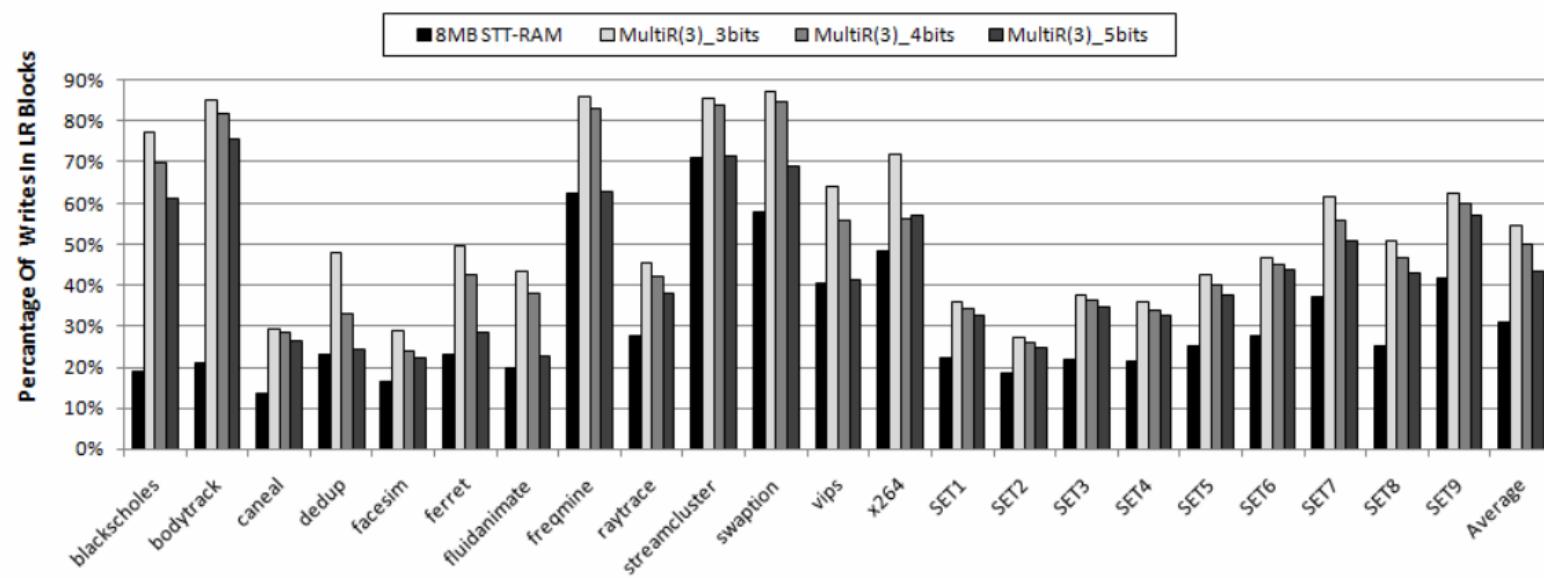


- 3 bit saturation counter



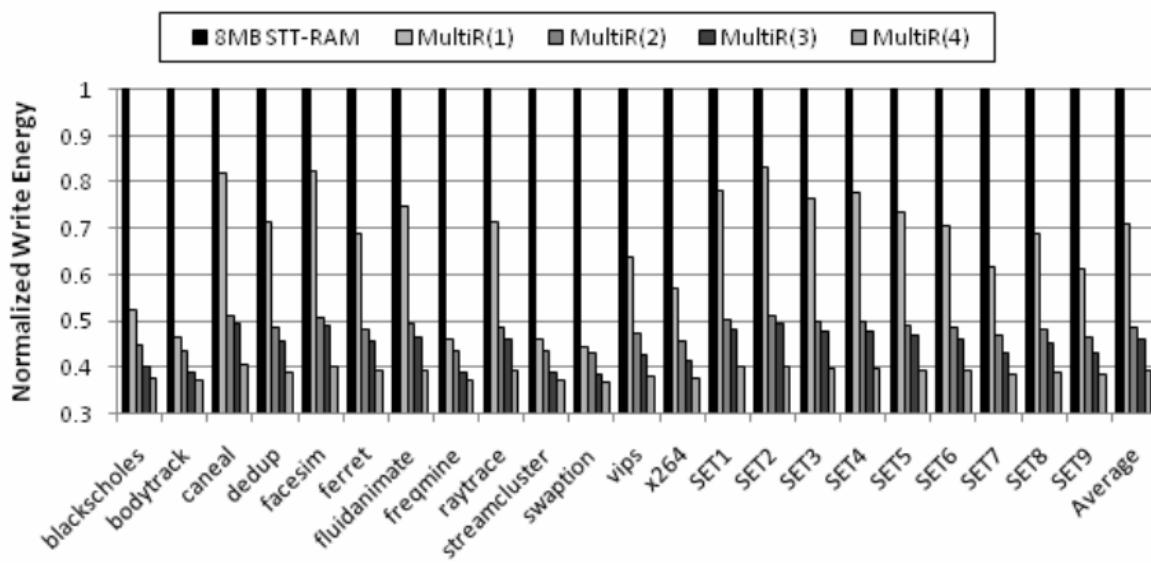
Sensitivity Analysis to Saturation Counter Bits

- 3 bit saturation counter

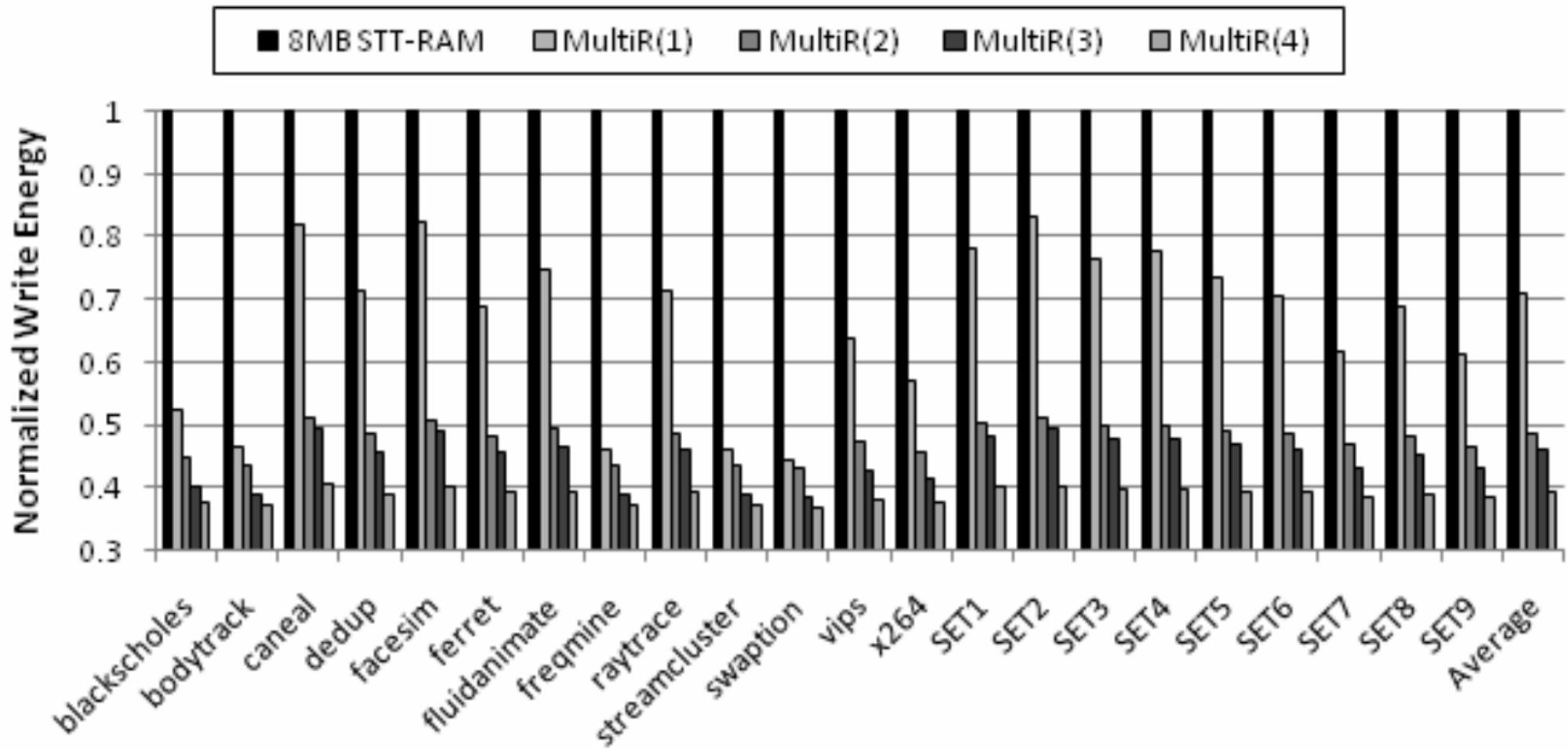


Write Dynamic Energy

- more than 55%

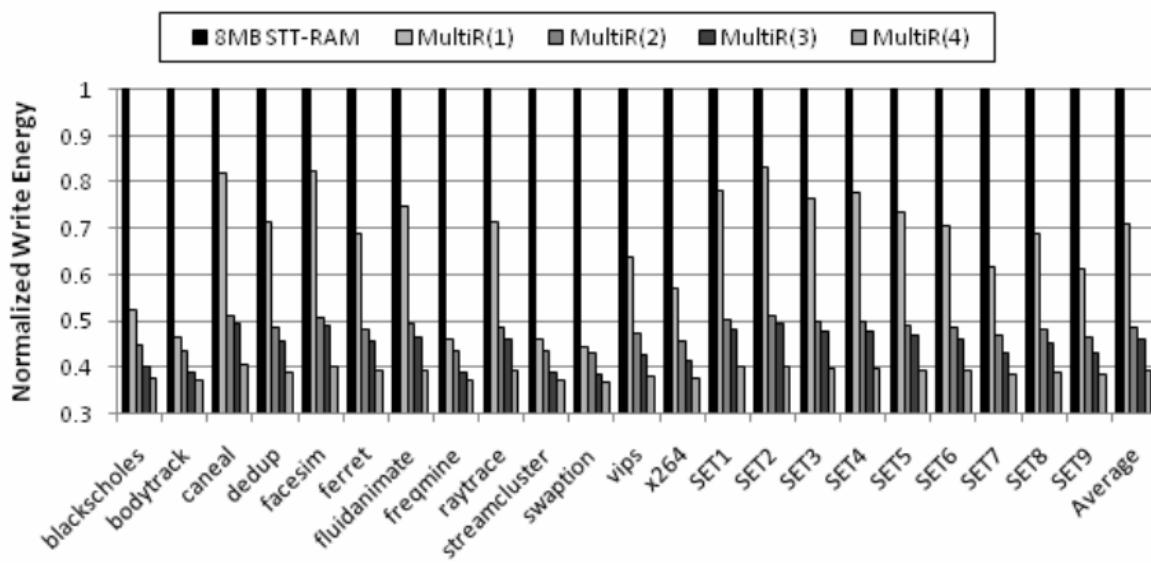


- more than 55%



Write Dynamic Energy

- more than 55%



Evaluation

Simulation Platform

- System-level simulation
 - Virtutech SIMICS full-system simulator
 - + Microblaze
 - Cycle accurate
 - + CACTI timing model rounded by system frequency
 - + CACTI energy model
 - Benchmarks: ROI of the SPEC CPU2006 + PARSEC-2 programs

Parameter	Description
PERIOD	Period of simulation. Default value is 100.
L1 Cache	Optimal size, 1024 bytes. Using this value L1 cache access latency is very low. Using lower memory, L1 cache access time increases.
L2 Cache	1024 MB ideal. Using smaller, 512 MB, memory latency increases.
L3 Cache	16 GB ideal. Using intermediate, 8 GB, size the access time is relatively poor.

21

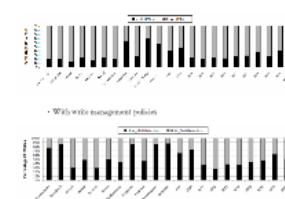
Design approaches

- 1.3 2MB SRAM replacement with 8MB STT-RAM with same area

Code	Op.	Execution Time	Dynamic Energy	Pulse Duration	Total Energy
I400	S	—	0.005 J/clock	—	0.005 J
STT-RAM	R	13 ns	0.007 J/clock	10 ns	0.007 J
	S	—	0.005 J/clock	—	0.005 J
I400	R	26 ns	0.005 J/clock	See	0.015 J
MRAM (Multi)	R	10 ns	0.007 J/clock	See	0.007 J
I400	R	10 ns	0.005 J/clock	10 ns	0.010 J
MRAM (Multi)	R	10 ns	0.007 J/clock	See	0.007 J
I400	R	10 ns	0.005 J/clock	10 ns	0.010 J
MRAM (Multi)	R	10 ns	0.007 J/clock	See	0.007 J
I400	R	10 ns	0.005 J/clock	10 ns	0.010 J
MRAM (Multi)	R	10 ns	0.007 J/clock	See	0.007 J
I400	R	10 ns	0.005 J/clock	10 ns	0.010 J
MRAM (Multi)	R	10 ns	0.007 J/clock	20 ns	0.014 J

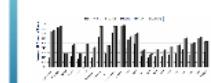
Write access distribution

- Without waste management policies



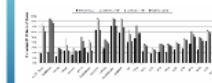
Write access distribution

- Industry Evolution
 - Technology Watch List that highlights STT-MRAM and its low cost vs. DRAM entry barriers (\$100 vs. \$10B)



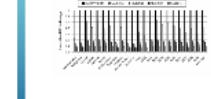
Sensitivity Analysis to Saturation Counter Bits

- 3 bit representation



Write Dynamic Env

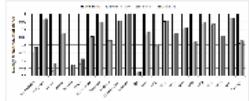
- more than 22%



Results

Average Memory Access Latency

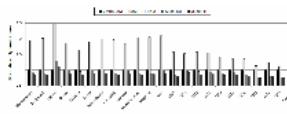
- AML = hit time + miss rate x miss penalty
 - AML (20%)



28

Dynamic Power

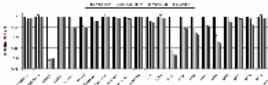
- 50% (STT-RAM)
 - 15% (SRAM)



1

Performance

- Performance (5%)



Total Power

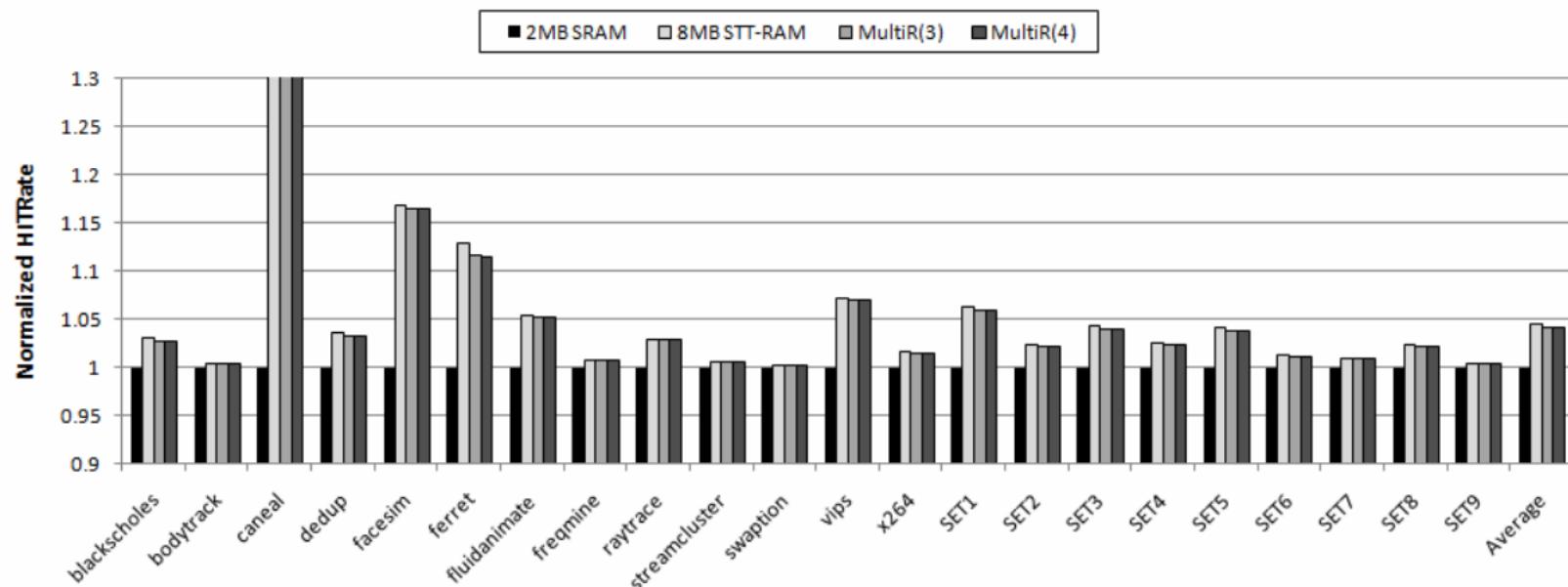
-

Results

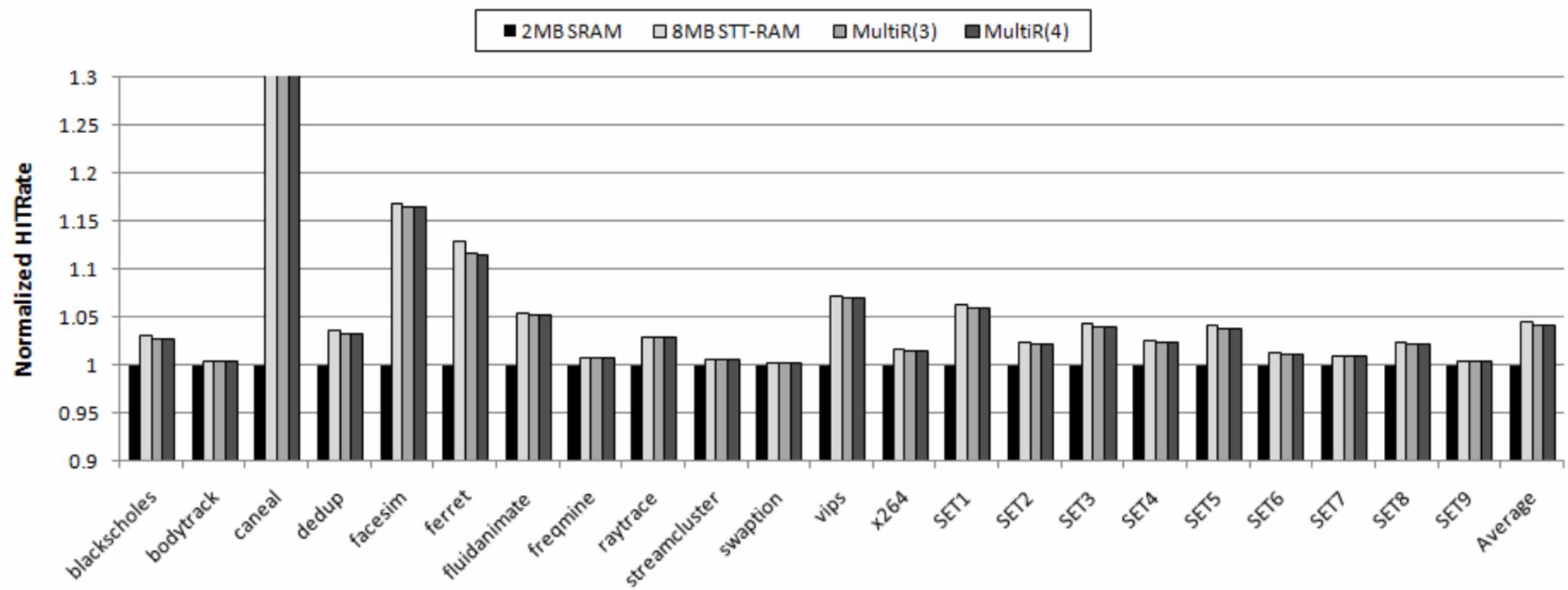
Hitrate

Hitrate

- hitrate (5%)

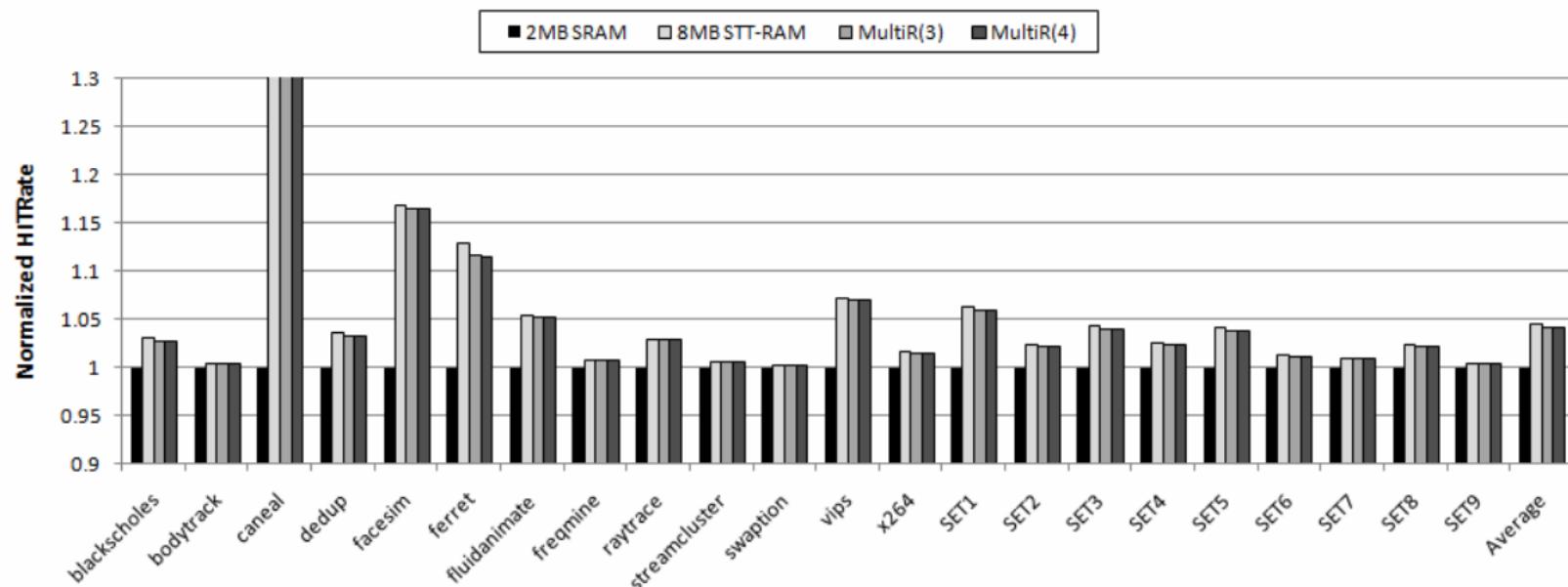


- hitrate (5%)



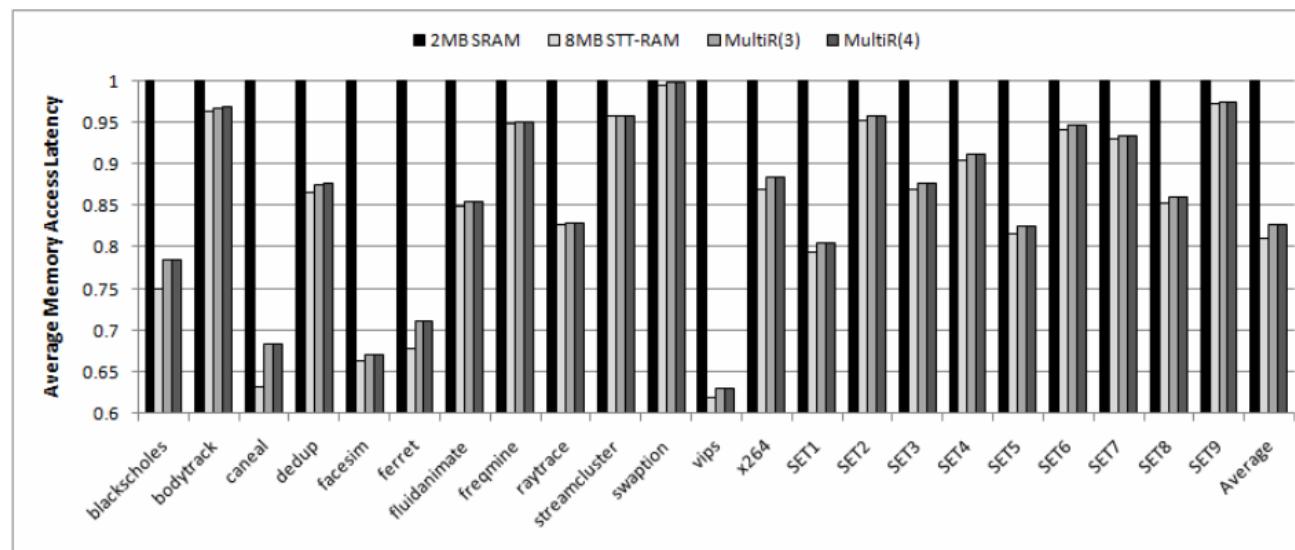
Hitrate

- hitrate (5%)



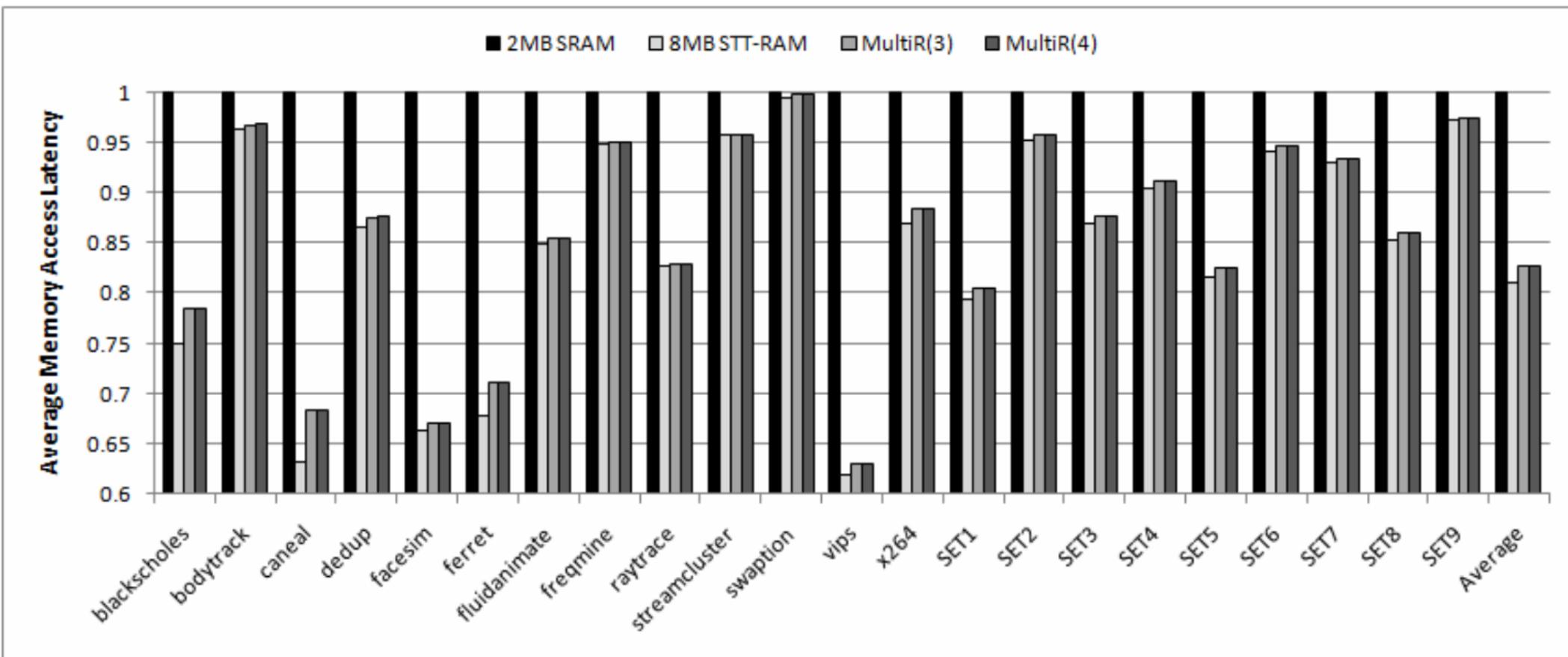
Average Memory Access Latency

- AML = hit time + miss rate x miss penalty
- AML (20%)



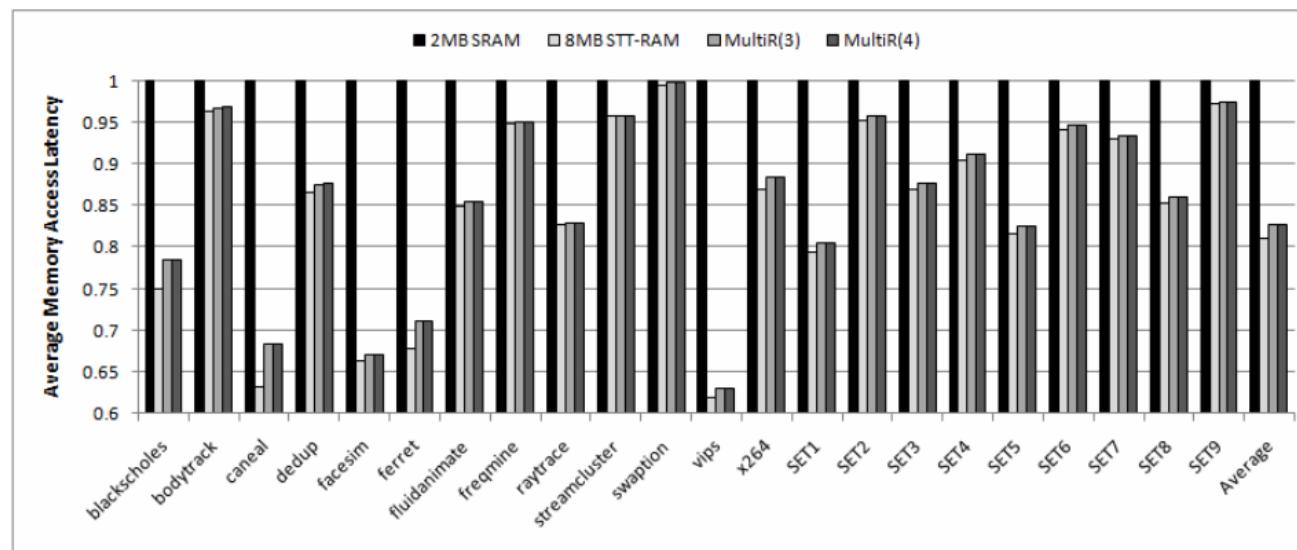
• $\text{AVIL} = \text{hit time} + \text{miss rate} \times \text{miss penalty}$

- AML (20%)



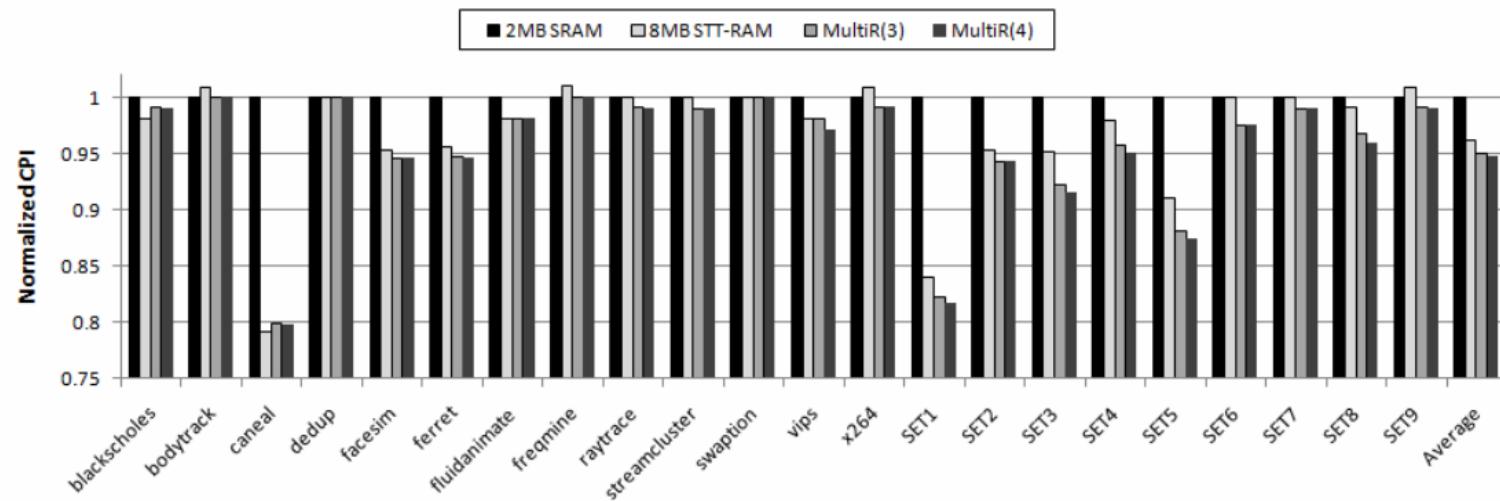
Average Memory Access Latency

- AML = hit time + miss rate x miss penalty
- AML (20%)

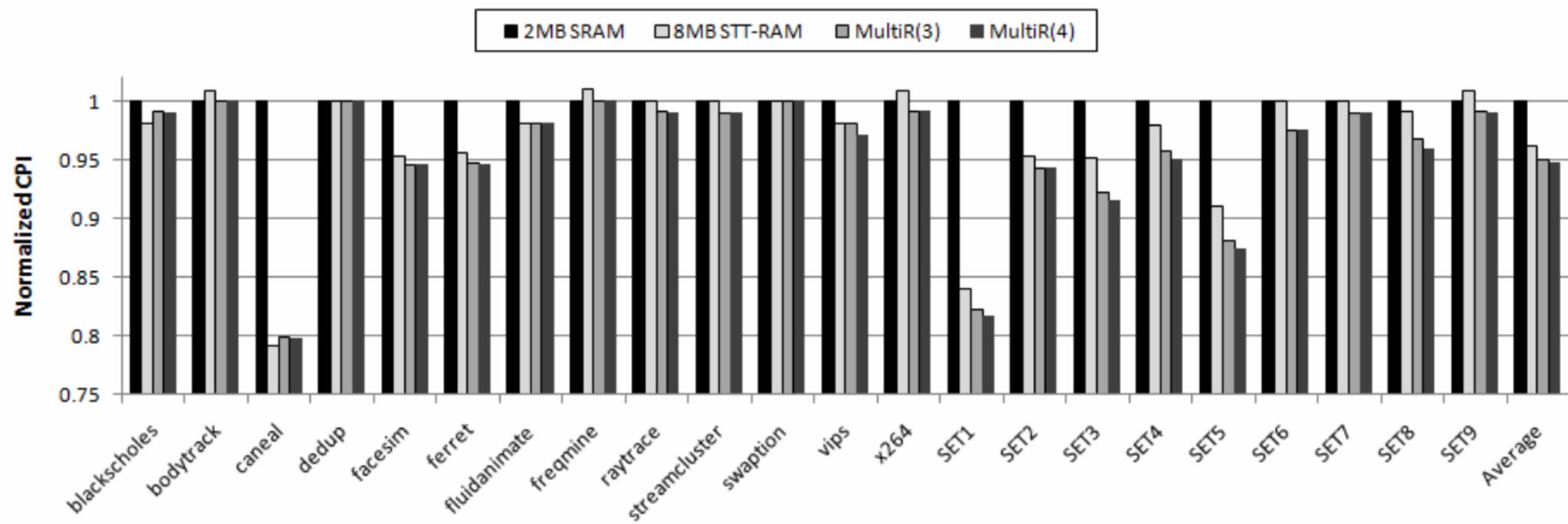


Performance

- Performance (5%)

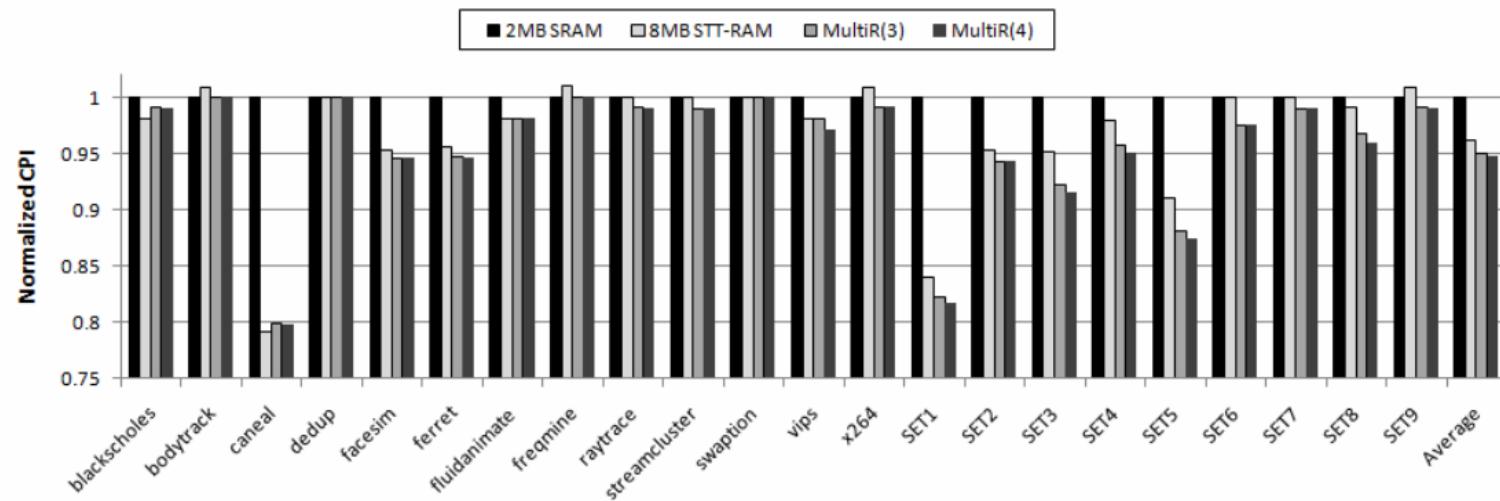


• Performance (5%)



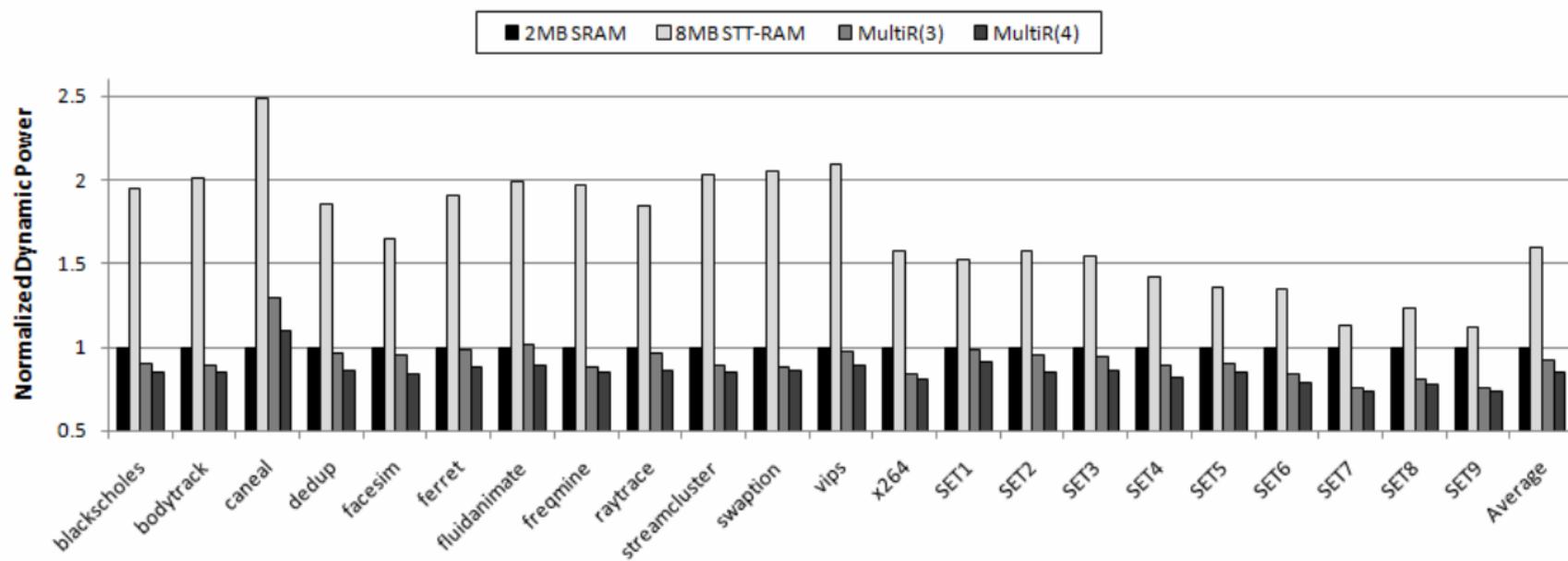
Performance

- Performance (5%)

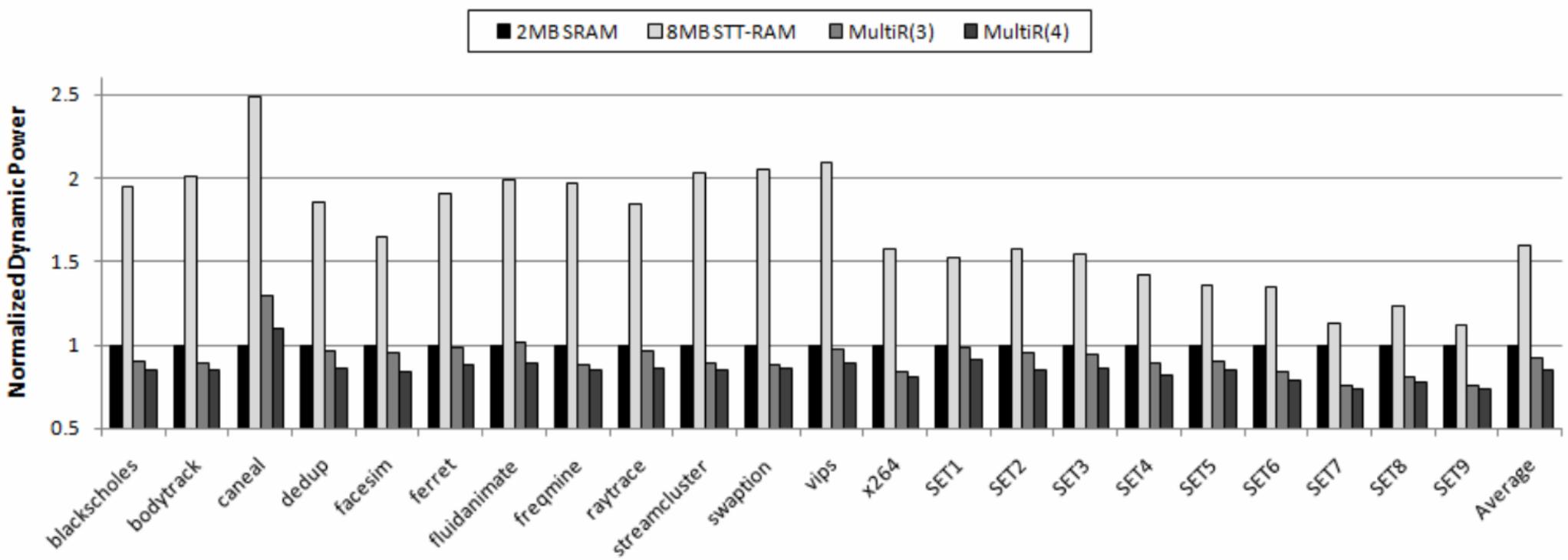


Dynamic Power

- 50% (STT-RAM)
- 15% (SRAM)



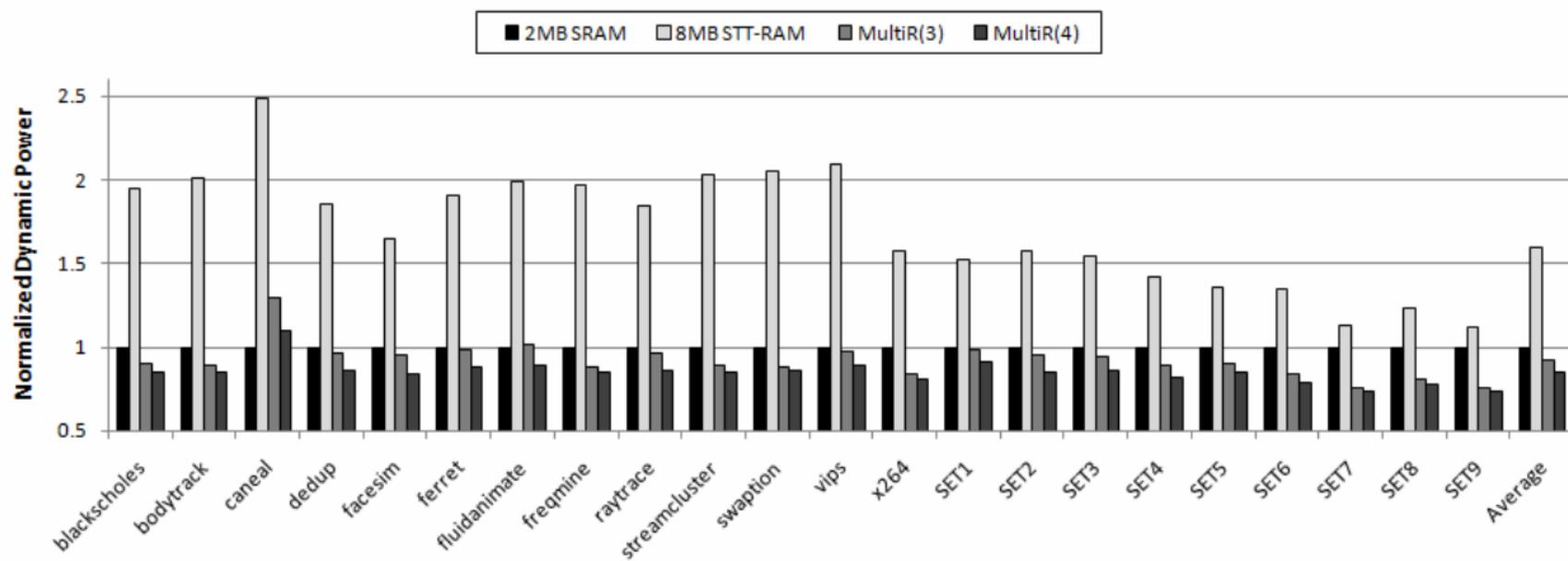
- 50% (STT-RAM)
- 15% (SRAM)



30

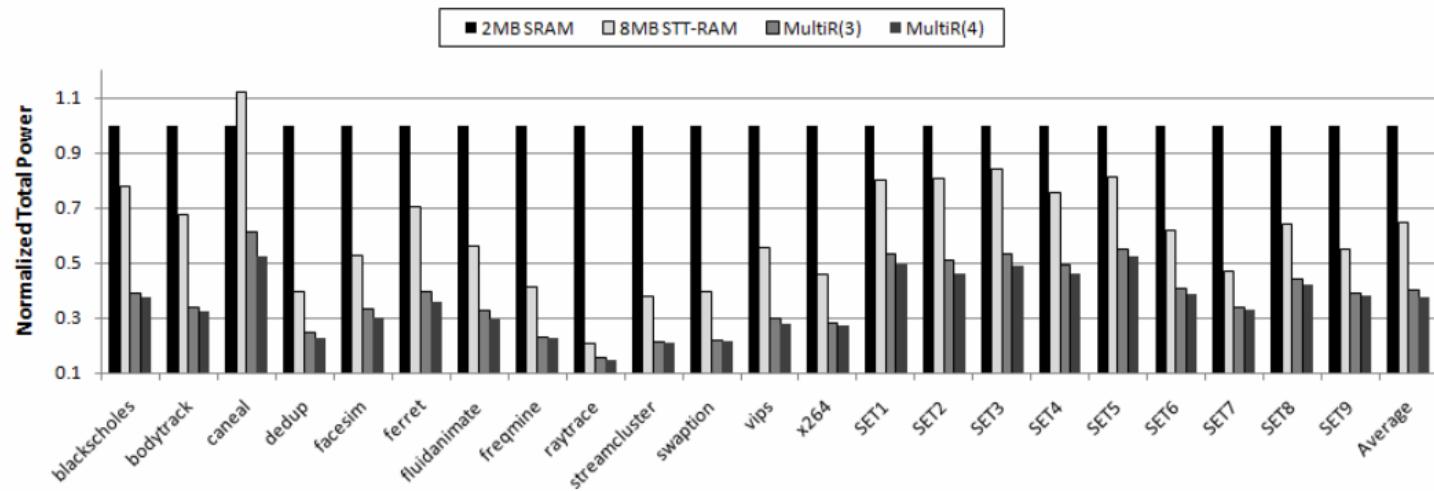
Dynamic Power

- 50% (STT-RAM)
- 15% (SRAM)



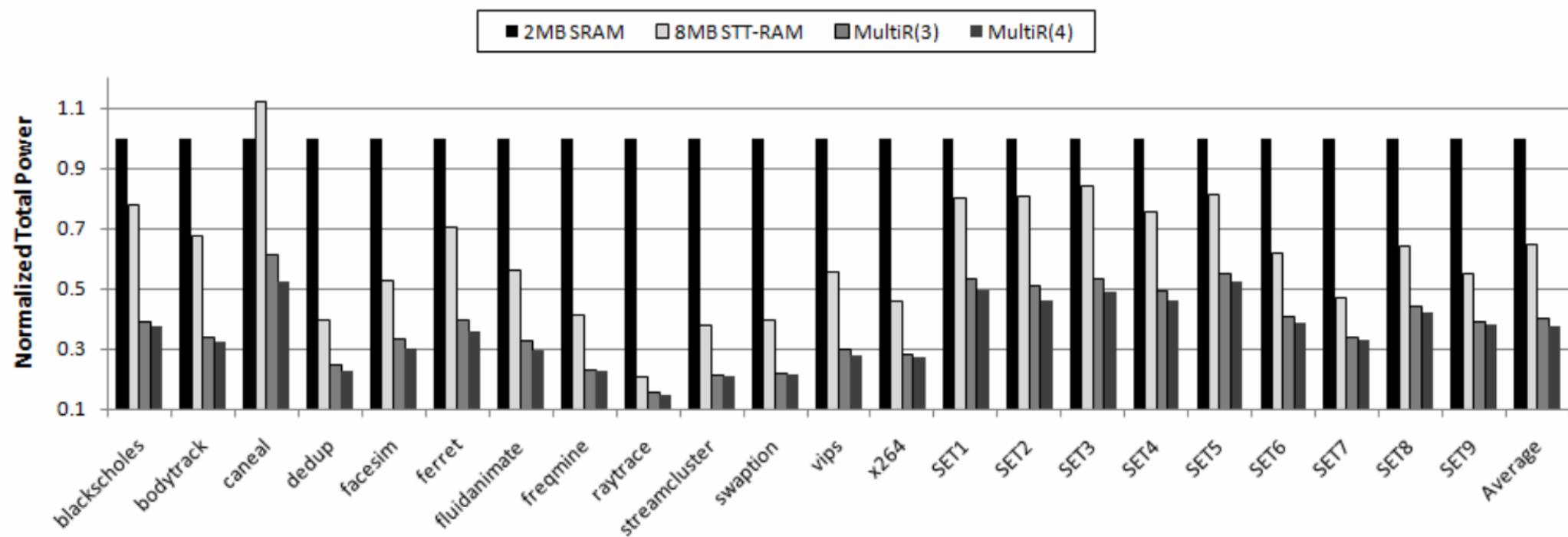
Total Power

- parsec (74%)
- spec (60%)



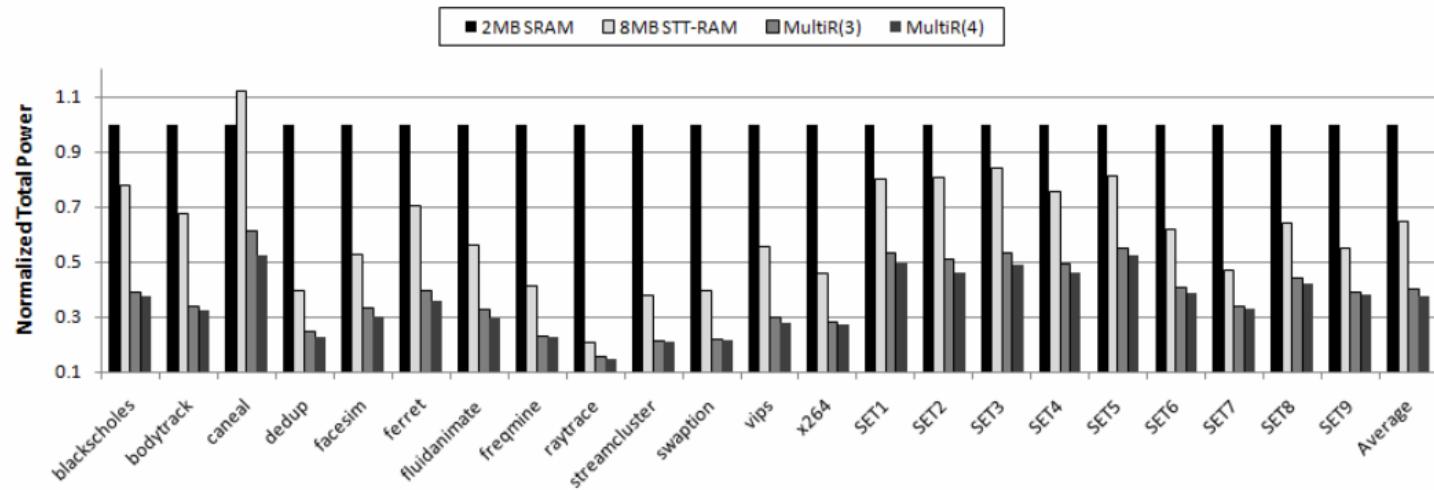
IUTAI I UWCI

- parsec (74%)
- spec (60%)



Total Power

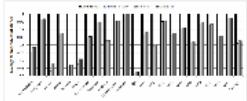
- parsec (74%)
- spec (60%)



Results

Average Memory Access Latency

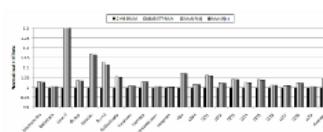
- AML = hit time + miss rate x miss penalty
- AML (20%)



28

Hitrate

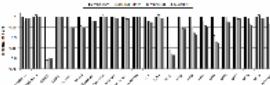
- hitrate (5%)



27

Performance

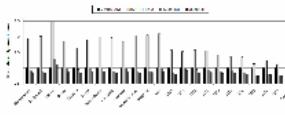
- Performance (5%)



29

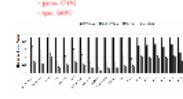
Dynamic Power

- 50% (STT-RAM)
- 15% (SRAM)



30

Total Power



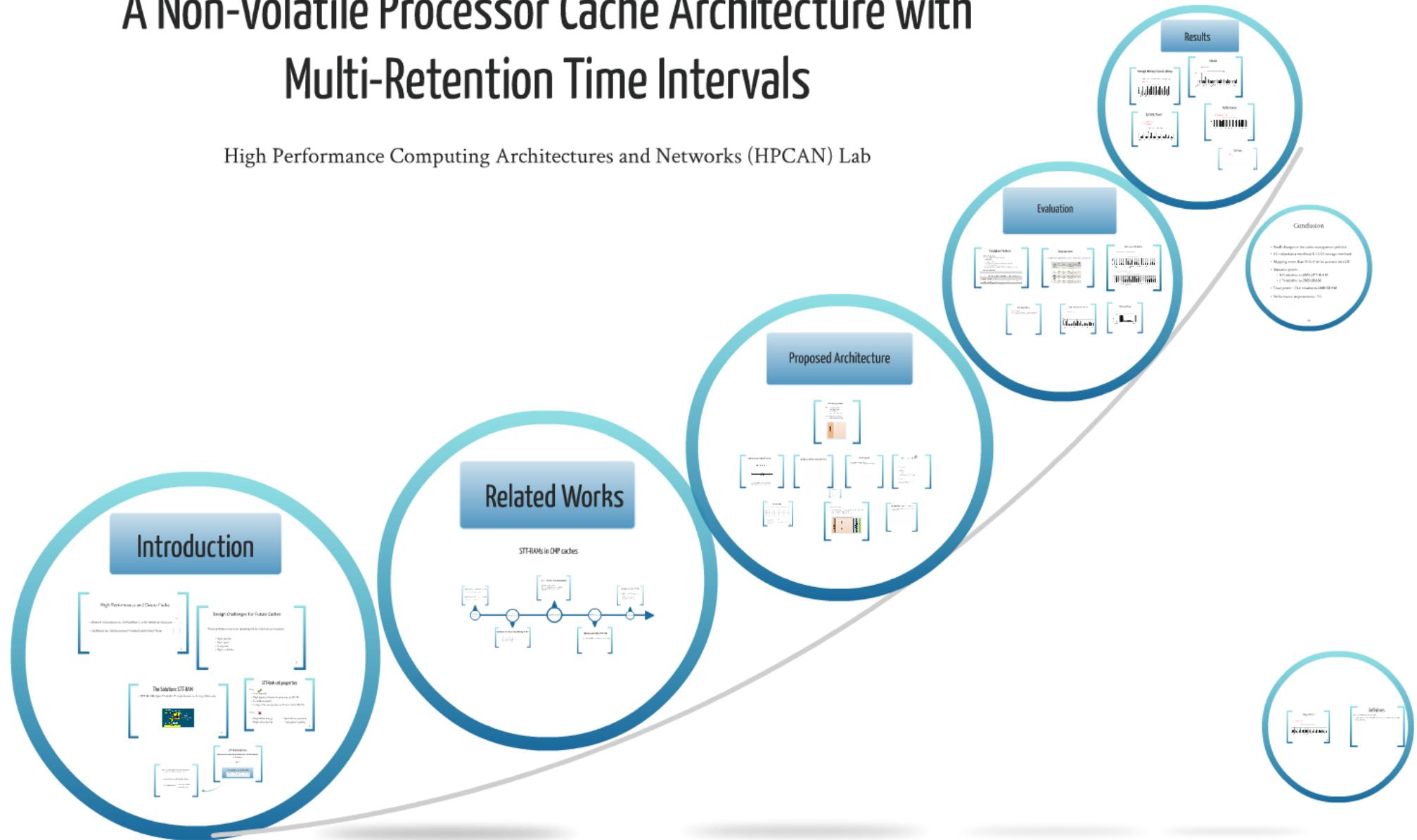
32

Conclusion

- Small changes to the cache management policies
- 1% redundancy overhead & 12.5% storage overhead
- Mapping more than 55% of write accesses into LR
- Dynamic power :
 - 50% relative to 8MB STT-RAM
 - 17% relative to 2MB SRAM
- Total power : 74% relative to 2MB SRAM
- Performance improvement : 5%

A Non-Volatile Processor Cache Architecture with Multi-Retention Time Intervals

High Performance Computing Architectures and Networks (HPCAN) Lab



Future works

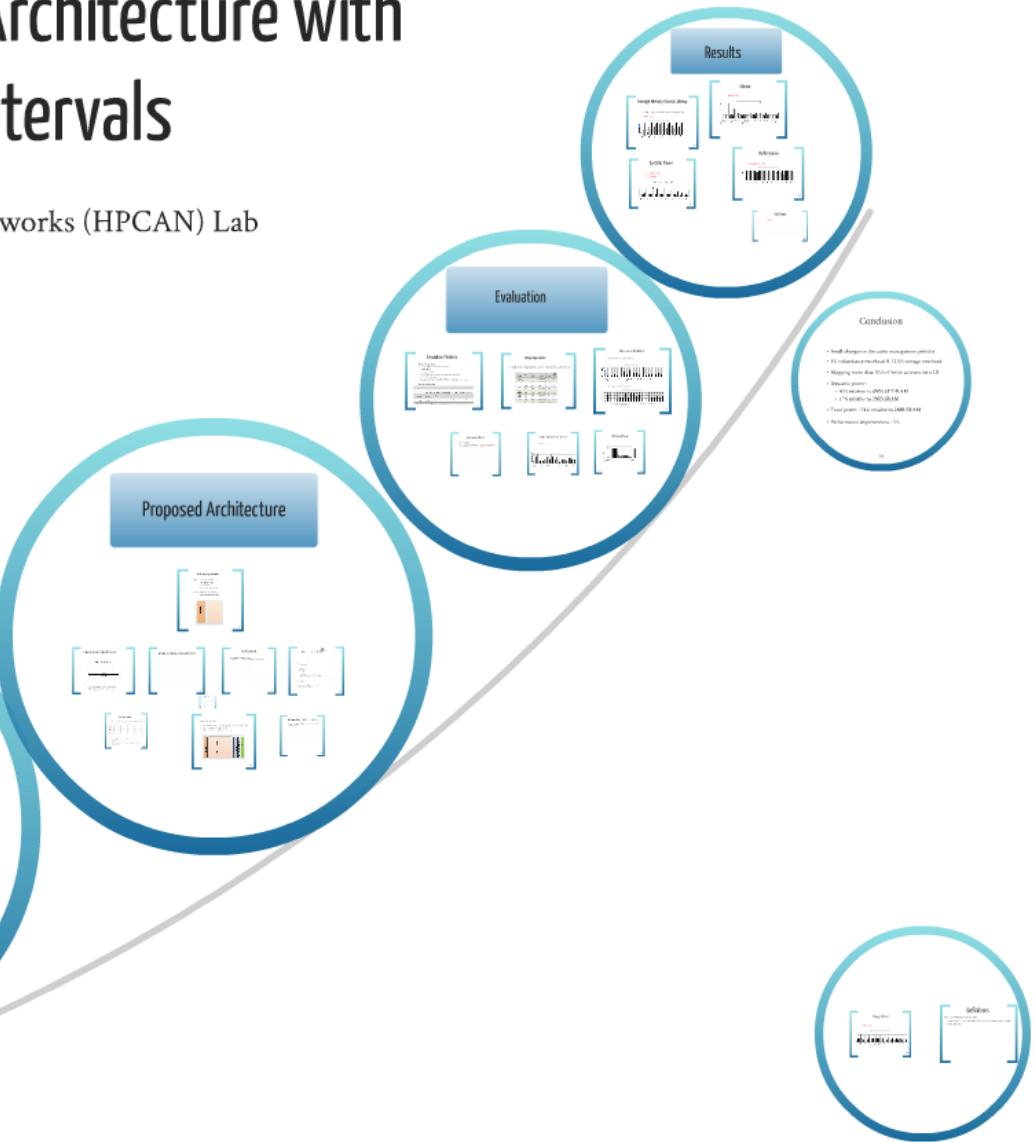
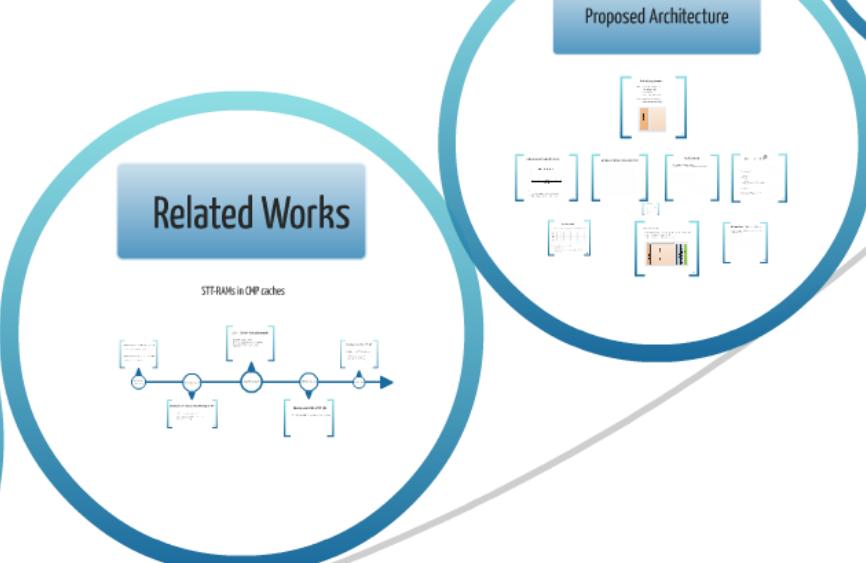
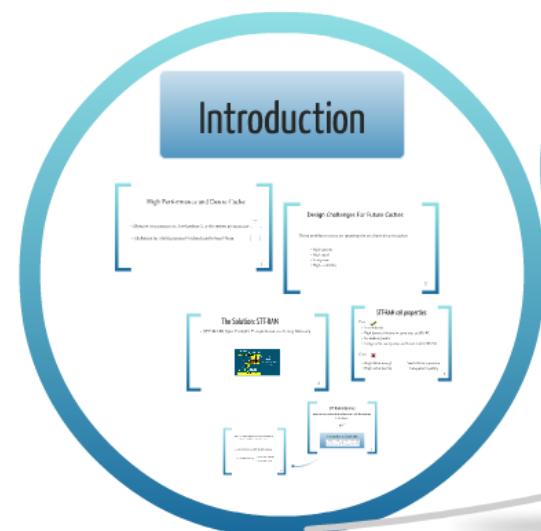
- Improve utilization of LR-region
 - Increase associativity of LR-region
- Different LR & HR array capacity
 - Increase LR ways in each set
 - Saving more write energy
- Different Refreshing mechanism for LR blocks
- Inter-set redirection:
 - increase number of hot blocks redirected to LR of destination set
- Considering endurance of STT-RAM LLC :
 - uniform distribution of write accesses

Thanks for your attention



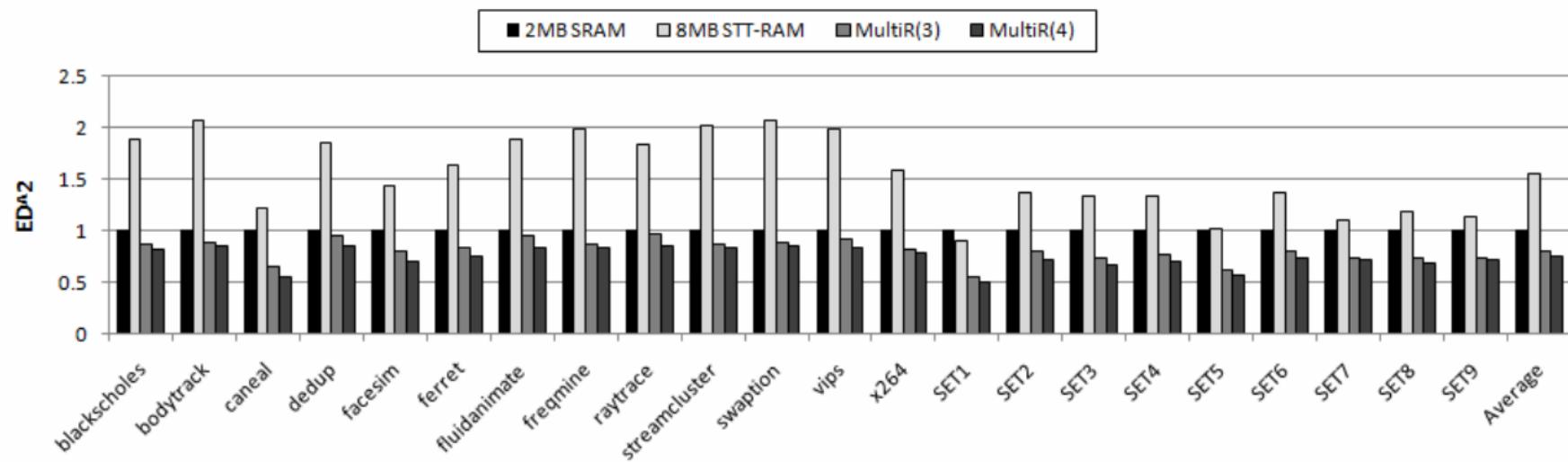
A Non-Volatile Processor Cache Architecture with Multi-Retention Time Intervals

High Performance Computing Architectures and Networks (HPCAN) Lab



Energy x Delay^{^2}

- ED^{^2} (25%)



- ED² (25%)

