

Mellow Writes: Extending Lifetime in Resistive Memories through Selective Slow Write Backs

ABSTRACT

Emerging resistive memory technologies, such as PCRAM and ReRAM, have been proposed as promising replacements for DRAM-based main memory, due to their better scalability, low standby power, and non-volatility. However, limited write endurance is a major drawback for such resistive memory technologies. Wear leveling (balancing the distribution of writes) and wear limiting (reducing the number of writes) have been proposed to mitigate this disadvantage, but both techniques only manage a fixed budget of writes to a memory system rather than *increase* the number available.

In this paper, we propose a new type of wear limiting technique, *Mellow Writes*, which reduces the wearout of individual writes rather than reducing the number of writes. *Mellow Writes* is based on the fact that slow writes performed with lower dissipated power can lead to longer endurance (and therefore longer lifetimes). For non-volatile memories, an N^1 to N^3 times endurance can be achieved if the write operation is slowed down by N times.

We present three microarchitectural mechanisms (*Bank-Aware Mellow Writes*, *Eager Mellow Writes*, and *Wear Quota*) that selectively perform slow writes to increase memory lifetime while minimizing performance impact. Assuming a factor N^2 advantage in cell endurance for a factor N slower write, our best Mellow Writes mechanism can achieve $2.58\times$ lifetime and $1.06\times$ performance of the baseline system. In addition, its performance is almost the same as a system aggressively optimized for performance (at the expense of endurance). Finally, *Wear Quota* guarantees a minimal lifetime (e.g., 8 years) by forcing more slow writes in presence of heavy workloads. We also perform sensitivity analysis on the endurance advantage factor for slow writes, from N^1 to N^3 , and find that our technique is still useful for factors as low as N^1 .

1. INTRODUCTION

DRAM technology scaling is fundamentally limited by its use of capacitance to store values, requiring energy wasting refreshes and losing values when power is removed. Emerging resistive memory technologies, such as resistive random access memory (ReRAM) and phase change memory have shown promise as DRAM replacements. Such emerging memory technologies have the advantage of high density, high scalability, non-volatility, and low standby power. They fall short, however, in one category: *write endurance*. Fatigued cells may fail to change state. This often results in

data errors, making write endurance a serious challenge for architecting resistive memory systems.

There are two common methods to combat the endurance limit of resistive memories:

- **Wear Leveling.** Most applications exhibit *non-uniform* write patterns, so the resistive memory cells in hotspot memory blocks will have a much shorter lifetime than others. Wear leveling evens out write patterns by remapping heavily written lines to less frequently written lines. The limit of wear leveling is that its lifetime improvement has an upper bound: the average lifetime of memory cells.
- **Wear Limiting.** Wear limiting attempts to reduce—rather than distribute—the amount of wear on the memory system. Some existing techniques are DRAM buffering [1] and Flip-N-Write [2]. All of these function by reducing the number of writes occurring to the resistive devices.

In this paper, we introduce a new technique to implement wear limiting on resistive memories. Instead of reducing the *number* of writes, we reduce the *impact* of some writes on the endurance by performing a slower write. In this paper, we will use an analytic model from [3] and ReRAM parameters, but our techniques are applicable to any systems with variable wear that is correlated with the speed of writes. In fact, our sensitivity analysis will show benefits for write speed to endurance relationships that vary from linear to cubic. The key, however, is to use slow writes strategically to avoid performance degradation.

We evaluate two *Mellow Writes* schemes: *Bank-Aware Mellow Writes* and *Eager Mellow Writes*. *Bank-Aware Mellow Writes* inspects the current set of write requests, sending slow writes to banks that have only one current write request. *Eager Mellow Writes* goes one step further, identifying useless dirty lines in the last level cache (LLC) to write back to banks with no requests. Experiments show that our proposed techniques preserve performance and significantly enhance lifetime—achieving almost the same performance as a system with the most aggressive speed-up techniques while achieving $2.58\times$ lifetime of a system with a default configuration without aggressive performance optimizations. In addition, we propose a *Wear Quota* to guarantee the minimal lifetime (e.g., 8 years in our experiment) of the ReRAM memory system. Finally, we find that while our scheme is influenced by the exponential relationship between latency

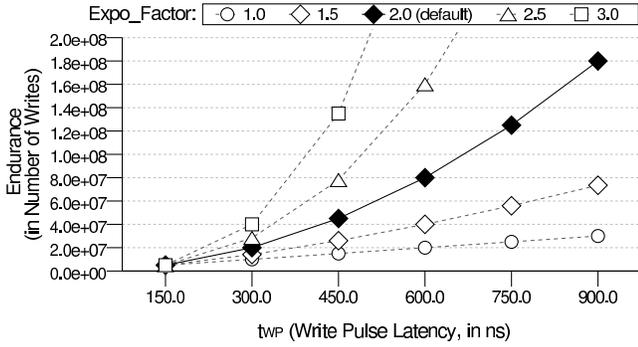


Figure 1: Trade-off between Write Latency and Endurance in our simulated resistive memory system (see Section 5). Without loss of generality, we use ReRAM technology with the baseline write latency of 150ns and the baseline endurance of $5 * 10^6$. For slow writes, we model the endurance based on Equation 2 (derived from paper [3]), and use five different *Expo_Factor*: 1.0, 1.5, 2.0, 2.5, 3.0. By default, we chose to model a quadratic write latency / lifetime tradeoff (*Expo_Factor* = 2.0).

and endurance we derive from [3], our scheme is still advantageous with a pessimistic linear model.

The rest of the paper is organized as follows. We first introduce the wear/latency trade-off in resistive memory in Section 2. Our motivation section (Section 3) shows the system-level performance impact of slow writes and the abundance of memory idle time. Section 4 describes our two *Mellow Writes* schemes and the *Wear Quota* scheme. We then present our methodology and results in Sections 5 and 6. We have a more detailed related work section in Section 7, followed by our conclusions in Section 8.

2. WRITE LATENCY/ENDURANCE TRADE-OFF

The relationship between write latency and endurance is intuitive and has been observed before [4, 5] for different kinds of non-volatile memories. In order to provide low latency in writes, it is typical to apply high power dissipation [6, 7]. Higher power dissipation accelerates failure mechanisms [8, 9, 5].

In this paper, we will use a recent analytic model [3] which makes two observations. First, switching speed is exponentially dependent on electric field and temperature for many types of nonvolatile memories including flash [10], phase change [11, 12], magnetoresistive and ferroelectric [13], and ReRAMs [14, 9] memories. Second, a high electric field combined with a high temperature exponentially increases the probability of creating new defects and/or filling existing deep traps in the dielectric [8] which is a primary source for limited endurance in majority of nonvolatile memories (e.g. those relying on electron-tunneling phenomena).

Combining these two observations, the analytic model [3] derives a polynomial relationship between endurance and write latency. Specifically, a cell’s endurance (in terms of total amount of writes) is given by:

$$Endurance \approx \left(\frac{t_{WP}}{t_0}\right)^{\frac{U_F}{U_S}-1} \quad (1)$$

where t_{WP} is write latency, t_0 is a device related constant, U_F is the activation energy for failure mechanism, and U_S is the activation energy of switching mechanism.

For non-volatile memories, practical values of U_S should be above 1eV [15]. Assuming practical values for U_F [8, 16], $\frac{U_F}{U_S}$ ranges from 2 to 4, so that, e.g., latency decrease is linearly, quadratically, and cubically proportional to endurance decrease for $\frac{U_F}{U_S} = 2$, $\frac{U_F}{U_S} = 3$ and $\frac{U_F}{U_S} = 4$, respectively. Thus our endurance model becomes:

$$Endurance \approx \left(\frac{t_{WP}}{t_0}\right)^{Expo_Factor} \quad (2)$$

where $1 \leq Expo_Factor \leq 3$.

As is shown in Figure 1, without loss of generality, we model a resistive memory with the baseline write latency of 150ns and the baseline endurance of $5 * 10^6$. Our baseline experiments model a quadratic write latency / lifetime trade-off (*Expo_Factor* = 2.0). This corresponds to $U_F \gtrsim 3eV$ which is representative of energy for creating a vacancy in many relevant metal oxide devices [17]. Our sensitivity experiments, however, model five different *Expo_Factor*: 1.0, 1.5, 2.0, 2.5, 3.0.

3. MOTIVATION

The motivation for our solution is based on two observations. First, using a single write latency cannot fulfill the performance and lifetime requirements for varying workloads. Second, for a system with a typical write latency, the memory banks are idle for much of the time.

In order to demonstrate the impact that different write latencies have on performance and lifetime, we run our baseline system (see Section 5 for details) with four different write latencies: normal write (1.0× latency) and slow writes (1.5×, 2.0× and 3.0× normal write latency). In addition, the use of write cancellations [18, 19] (in the presence of a read to the same bank) may also influence both performance and lifetime, so we perform simulations with and without write cancellation.

Figure 2 shows the results. We make the following important observations:

- Short latency (e.g., 1.0-1.5× normal writes) leads to unreasonably short lifetimes for some benchmarks (e.g., `lbm`, `leslie3d`, etc.).
- Although slow writes provide longer lifetime, they can considerably degrade the overall system performance (`stream`: 63.8% degradation at 3.0×, 30.1% at 1.5×. Therefore, slow writes must be used judiciously.
- Write cancellation is no silver bullet for slow writes. It helps by allowing reads to complete more quickly (`milc`, `mcf`), but can degrade performance by putting pressure on write queues, leading to more of the expensive write drain operations (`hammer`, `bwaves`). Write cancellation also comes at a penalty to memory lifetime due to the multiple write attempts. Therefore, write cancellation is only part of the solution to reduce the performance penalty of slow writes.

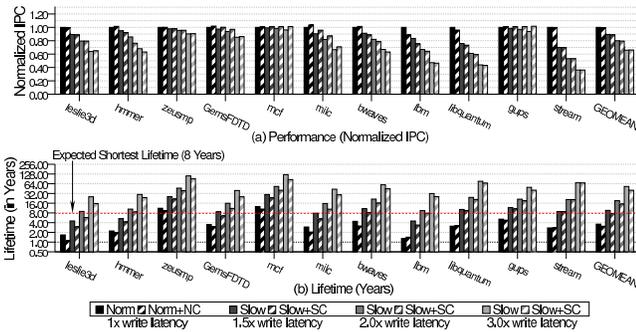


Figure 2: Normalized IPC and Lifetime (in years) of systems with normal writes and $1.5\times$ – $3.0\times$ slow writes.

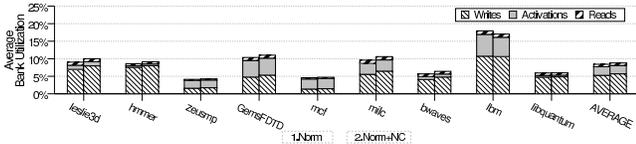


Figure 3: Average bank utilization of systems with normal writes. The utilization of a bank refers to the percentage of the time when corresponding bank is busy.

In summary, to get higher performance on the premise of guaranteeing minimal lifetime (in our case, 8 years), different applications favor different write latencies and write policies. For example, *les1ie3d* favors $2\times$ latency writes without write cancellation; *gups* is suitable for $1.5\times$ latency writes with write cancellation; and *zeusmp* likes $1\times$ latency writes without write cancellation. As a result, it is hard, if not impossible, to find a fixed write latency and a fixed write policy to fit all the applications.

Therefore, it is necessary to use adaptive write schemes in which a system can adaptively use fast and slow writes. The former ones improve the performance, and the latter ones extend the lifetime. Figure 3 shows the opportunity—in a system with fast writes (i.e., with $1.0\times$ latency), the utilization of memory banks is not particularly high. Therefore, there are ample opportunities to selectively write dirty data using slower writes (e.g., with $3.0\times$ latency) when their corresponding banks are not busy.

Our goal is to use slow writes at times when it is least likely to lead to performance degradation. Intuitively, we will be using the LLC as a large buffer from which we can find proper cachelines to be eagerly and slowly written back during periods of memory idle time. Predicting which cachelines are proper to do so, and when to perform the eager slow writes, is the challenge we face.

4. MELLOW WRITES

In this section, we discuss how to adaptively use fast and slow writes. For hardware simplicity consideration, in this paper we just adaptively use two different kinds of writes: normal writes with $1\times$ latency, and slow writes with $3\times$ latency.

We introduce our two *Mellow Writes* schemes that selectively perform slow writes when the memory system is rel-

atively less busy. Both schemes depend on idle memory cycles to perform such writes. In order to protect against memory-intensive workloads that do not have enough idle times, we also introduce a *Wear Quota* scheme to provide guaranteed lifetimes (at the cost of performance).

4.1 Bank-Aware Mellow Writes

Bank-Aware Mellow Writes scheme exploits the fact that a program could have many writes and still have idle time in a particular bank. Therefore, we make decisions at the bank granularity. A write request can be issued as a slow write only as long as there are no other operations (reads or writes) queued for the same bank. Figures 4 and 5 illustrate this scheme.

Figure 4 shows a situation with several read and write requests. For Bank 1, there is a single write request and no read requests. Therefore, the write request for Bank 1 can be issued as a slow write.

Figure 5 shows a situation with two write requests for Bank 0 and no read requests for the same bank. In this case, the next write request for Bank 0 will be issued as normal write, since there is another write request waiting. This is to reduce the chances that the write queue will fill, triggering an expensive write drain.

Figures 4 and 5 also illustrate that no write requests can be issued to Banks 2 and 3, since there are read requests for the banks waiting. Read requests have higher priority than write requests.

One advantage of the *Bank-Aware Mellow Writes* technique is that it requires minimal changes to the memory controller. The only modifications are a mechanism to detect bank conflict in the read and write queues and implementation of the slow write technique. However, its glaring drawback is that, when there are no writes for a bank in the write queue, it is unable to take advantage of the bank idle time.

4.2 Eager Mellow Writes

In order to take advantage of bank idle time in the presence of no write requests, we introduce *Eager Mellow Writes*. This allows the cache to eagerly write back some dirty data (that are predicted to not be used again) when the memory is not busy. In a nutshell, these items are placed into a third queue (*Eager Mellow Queue*). Items in *Eager Mellow Queue* have the least priority, can never trigger a write drain, and are only performed when there are no normal read or write requests to that bank.

Figure 6 shows a high level view of *Eager Mellow Writes* scheme. In this case, only slow write requests for Bank 2 can be issued from *Eager Mellow Queue*. Requests for other banks cannot be issued from *Eager Mellow Queue* because there are outstanding requests for these banks in the write and/or read queues.

This scheme requires two changes. First, the LLC needs to identify what items to use as candidates for *Eager Mellow Writes*. These are sent to the memory controller. Second, the memory controller needs an additional *Eager Mellow Queue* to hold them. In the rest of this subsection, we will discuss in detail the design of *Eager Mellow Writes* technique.

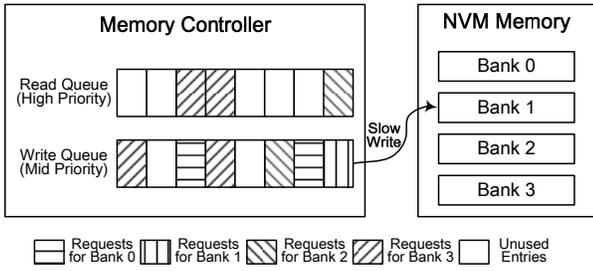


Figure 4: Situation when *Bank-Aware Mellow Writes* scheme issues a slow write.

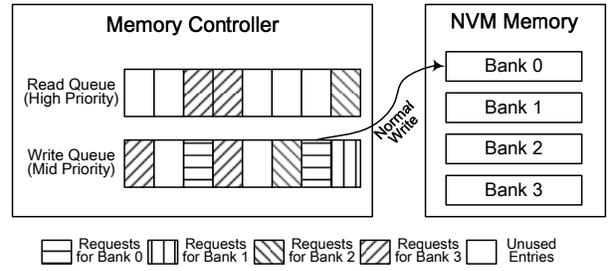


Figure 5: Situation when *Bank-Aware Mellow Writes* scheme issues a normal write.

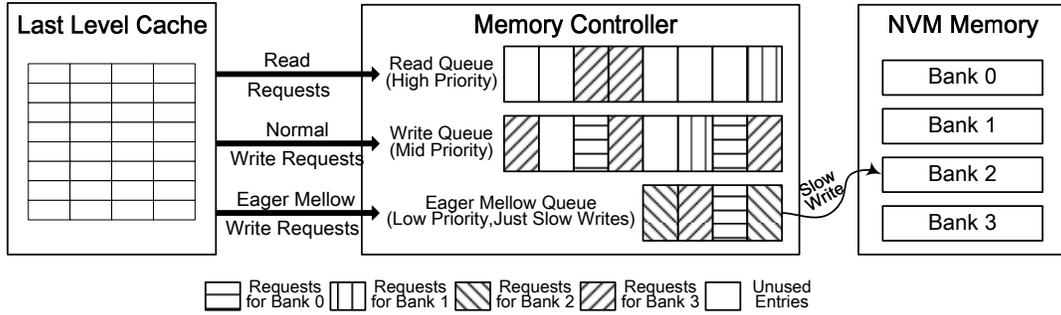


Figure 6: A high level view of a processor and a memory controller using *Eager Mellow Writes* technique.

4.2.1 Identifying Eager Mellow Writes

Any cycle the LLC is idle and the *Eager Mellow Queue* is not full, the LLC has a chance to choose an item to be placed in the *Eager Mellow Queue*. The goal is to find dirty lines in the cache that are unlikely to be modified again before they get evicted from the cache. If a line gets modified again before being evicted, then the write was wasted, reducing, rather than increasing, the lifetime.

Our scheme is as follows. The LLC randomly chooses a cache set. Within that set, it finds the dirty lines that are unlikely to be used again. If present, it issues an *Eager Mellow Write* of the least likely line to be used again. That line is marked clean, *not evicted*, from the cache.

The key design issue—how to find the useless dirty cache line, in a simple, energy-efficient, low-storage way? Given the stack property [20] of LRU replacement policy, and inspired by Qureshi *et al.*'s work [21][22], we propose a simple but effective scheme for an N-entry LRU stack (as shown in Figure 7):

- In the LLC controller, add a hit counter for every LRU stack position. (Note: This is a single counter for the same LRU stack position across all sets, not per set). For an N-way associative LLC set, the *Most Recently Used* block is in LRU position 0, while *Least Recently Used* block is in LRU position (N-1). In addition, add a single miss counter to record the number of missed requests.
- On every LLC request (read or write), based on the hit stack position, increment the corresponding hit counter on a hit or miss counter on mis.
- After every T_{sample} period, we check the count of all the counters, identify the eager LRU position. This is

the position such that the sum of the hits in eager LRU position through n-1 is less than $THRESHOLD_RATIO$ of the requests. Since these higher LRU positions contribute so few hits, they are marked as *useless* until the next check and are subjected to *Eager Mellow Writes*. After we set the new *useless* LRU stack positions, we reset the hit counters and miss counter to 0, and restart the new round of profiling for the next T_{sample} period. In our experiment, T_{sample} is 500000 ns, and $THRESHOLD_RATIO$ is $\frac{1}{32}$. As is shown in our motivational example of Figure 7, LRU positions 3—7 accumulate less than $\frac{1}{32}$ of the total LLC requests and therefore considered as *useless* and subjected to *Eager Mellow Writes*.

4.2.2 Performing Eager Mellow Writes

We add a third category of memory accesses to the memory controller. *Eager Mellow Writes* from the LLC are placed into the *Eager Mellow Queue*, which can only issue slow writes to banks. This queue has the lowest priority, no write drain operations, and is only issued when there are no same-bank requests in either of the read queue and the normal write queue. Also, to reduce the hardware overhead as well as reduce the number of *Eager Mellow Writes*, we use a relatively small *Eager Mellow Queue*—in our simulation, it has only 16 entries, whereas the read and write queues each have 32 entries.

4.3 Wear Quota

As will be shown in Section 6, although *Mellow Writes* greatly improve the lifetime of a resistive main memory system without noticeable performance penalty, with a memory-intensive workload, the lifetime can still fall below an acceptable threshold (e.g., 8 years). Here, we introduce an

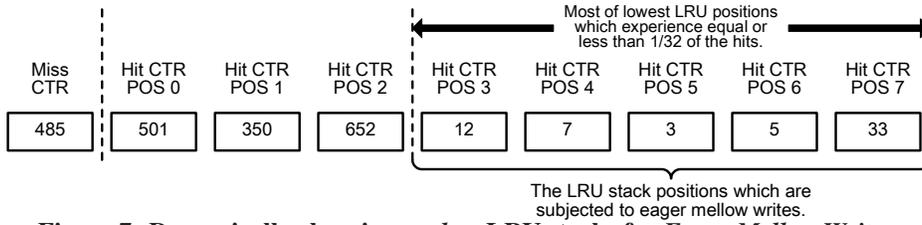


Figure 7: Dynamically choosing useless LRU stacks for *Eager Mellow Writes*.

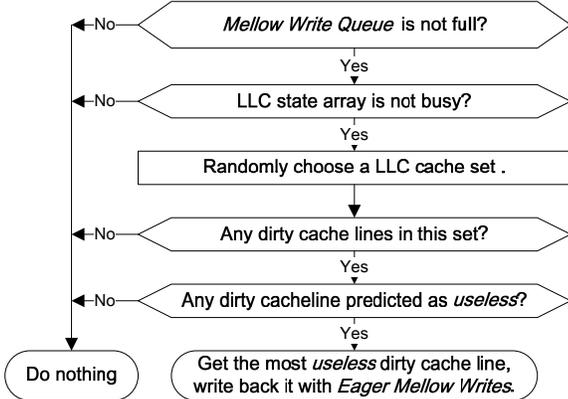


Figure 8: Workflow of *Eager Mellow Writes* on LLC side.

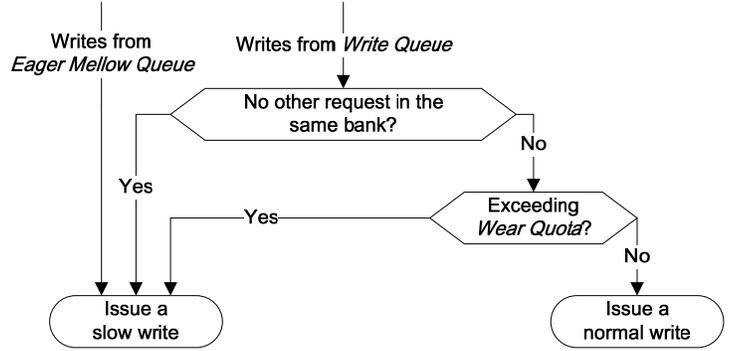


Figure 9: How a memory controller using *Bank Aware Mellow Writes*, *Eager Mellow Writes*, and *Wear Quota* decides whether to use normal writes or slow writes.

other scheme, *Wear Quota*, to guarantee the lifetime for memory write intensive applications. The notion of *Wear Quota* is straight-forward: We divide the execution period into multiple sample periods (e.g., 500000ns per period in our experiment). If the accumulative wear of all previous periods surpasses their corresponding wear threshold, only slow writes can be issued in current period.

In our resistive memory system, the write operation is at a block granularity (64 bytes), and the endurance of each resistive memory block is $Endur_{blk}$ (in terms of normal writes). If we want the lifetime of a block to be at least $T_{lifetime}$, we can ensure its average wear during every period of T_{sample} to be at most $WearBound_{blk}$:

$$WearBound_{blk} = Endur_{blk} * \frac{T_{sample}}{T_{lifetime}}$$

When referring to the lifetime of a memory bank, its wear in each sample period could be on average at most $WearBound_{bank}$:

$$WearBound_{bank} = BlkNum_{bank} * WearBound_{blk} * Ratio_{quota}$$

where $BlkNum_{bank}$ is the number of memory blocks of the memory bank, and $Ratio_{quota}$ is a number equal or smaller than 1.0. Ideally, $Ratio_{quota}$ should be 1.0. However, in our experiment (and also necessarily in real system), we guarantee even write distribution by using Start-Gap wear leveling technique. Since Start-Gap may introduce slightly extra wear, we conservatively set $Ratio_{quota}$ to be 0.9.

In the beginning of each period, the memory controller first calculates the value of $ExceedQuota$ of each memory bank:

$$ExceedQuota = \sum Wear_{bank} - WearBound_{bank} * Num_{previous_periods}$$

where $\sum Wear_{bank}$ is total amount of wear placed on the memory bank, and $(WearBound_{bank} * Num_{previous_periods})$ is the *Wear Quota* of all previous periods. If $ExceedQuota$ is larger than 0, it means total wear of the previous periods exceeds corresponding quota. **In this case, to reduce the average amount of wear per period, this memory bank can only performance slow writes in the coming period.**

4.4 Put It All Together

Our three proposed schemes (*Bank-Aware Mellow Writes*, *Eager Mellow Writes* and *Wear Quota*) work together to deliver high performance most of the time, but still guarantee a certain memory lifetime. Figure 9 shows how a memory controller with these three schemes decides when to issue a slow write: *For each bank*, look for a write to perform.

- If there is a single request in the *Write Queue*, issue a slow write.
- If there are multiple requests, but the *Wear Quota* is exceeded, issue a slow write.
- If there are multiple requests and the *Wear Quota* is not exceeded, issue a normal write.
- If there are no requests in the *Write Queue*, and there is a request in the *Eager Mellow Queue*, issue a slow write.

4.5 Overhead Discussion

In this subsection we will discuss the hardware and energy overhead of proposals.

- **Additional Voltage Supply.** Our proposed scheme requires a lower voltage supply to enable the slow write operation with smaller dissipated power. Since resistive memory circuits typically already require more than

one voltage supply (e.g., write and read operations typically need different voltages), there is a negligible design overhead.

- Storage Overhead.** The LLC, for *Eager Mellow Writes*, requires some additional storage—a cycle counter, a miss counter and number of LLC associativity hit counters, where each counter is $\lceil \log_2 \frac{T_{sample}}{T_{proc_clk}} \rceil$ bits. In our experiments, the LLC has an associativity of 16, T_{sample} is 500000ns, and T_{proc_clk} (processor clock period) is 0.5 ns. The total additional storage in LLC is $\lceil \log_2 \frac{T_{sample}}{T_{proc_clk}} \rceil * (1 + 1 + ASSOC_{LLC}) = 20 * 18 = 360$ bits. The memory controller also request additional storage. *Eager Mellow Writes* requires a 16-entry queue for each memory channel, and *Write Quota* requires three registers (64-bit each) in each bank to record the total number of normal writes, slow writes and periods to the corresponding memory bank.
- Energy Overhead.** Additional energy is used in the LLC and memory. When finding useless dirty cache lines, the LLC state array is the only RAM needed to be accessed; LLC tag/data RAMs will be accessed only when a useless dirty cache line is found and needed to be issued as a mellow eager writeback. Our technique introduces extra energy consumption in memory because (1) *Eager Mellow Writes* and *Write Cancellation* generate extra number of writes, and (2) a slow write consumes more energy than a normal write. However, in Section 6 we will quantitatively show that the additional memory-side energy consumption is moderate compared with whole system energy.

Overall, compared with the lifetime benefit, our design requires minimal hardware overhead and a moderate increase in main memory energy consumption.

5. METHODOLOGY

We use the gem5 simulator [24] with NVMain [25], a timing-accurate main memory simulator for non-volatile memory technologies. Table 1 provides the processor and cache details, and Table 2 provides the memory system details. For a given workload, we assume the system will cyclically execute the same execution pattern. Then the lifetime is calculated as how much time it takes until one cell in the memory system reaches its wear limit.

Without loss of generality, we model ReRAM technologies [26] which have a representative *Expo_Factor* of 2.0. ReRAM represents a wide range of technologies, with their write latency ranging from few nanoseconds [27] to millisecond scale [28], and their endurance ranging from few hundreds [29] to 10^{12} scale [30]. Here we consider representative memory-grade ReRAM devices with 150ns normal write latency and $5 * 10^6$ normal write endurance. Note that, recently Intel and Micron announced their partnership to create Xpoint Memory [31], while HP and SanDisk teamed up to build their SCM (Storage Class Memory) [32]. Both technologies are based on non-volatile memory (SCM based on ReRAM, but Xpoint Memory’s exact technology unknown),

Table 1: Processor Simulation Parameters

Freq.	2GHz
Core	OoO, Alpha ISA, 8-issue, 64-byte cacheline
# of Cores	For single-program: single core. For multi-program: four cores.
L1\$	split 32KB I/D-caches per core, 4-way, 2-cycle hit latency, 8-MSHR
L2\$	256KB per core, 8-way, 12-cycle hit latency, 12-MSHR
L3\$ (LLC)	16-way, 35-cycle hit latency, useless threshold is $\frac{1}{32}$, profiling period is 500,000ns. For single-program: 2MB, 32-MSHR For multi-program: 8MB, 80-MSHR

Table 2: Main Memory System Simulation Parameters

Basics	400 MHz, 64-bit bus width, using ReRAM, using Start-Gap wear-leveling scheme [23] in bank granularity, write-through (writes bypass row buffers), 1KB row buffer, open page policy, tFAW=50ns
# of Channels	For single-program: 1 For multi-program: 4
# of Banks per Channel	Three options: –4 banks, distributed in 1 ranks –8 banks, distributed in 2 ranks –16 banks, distributed in 4 ranks (default)
Read Queue	32 entries per channel, highest priority
Write Queue	32 entries per channel, middle priority, write drain threshold: 16 (low), 32 (high)
Eager Mellow Write Queue	16 entries per channel, lowest priority, no write drain, slow writes
Wear Quota Parameters	Expected Lifetime: 8 Years <i>Wear Quota</i> sample period: 500,000ns <i>Wear Quota</i> threshold ratio ($Ratio_{quota}$ in Section 4): 0.90
Row Size	16KB
tRCD	48 cycles (120 ns)
tWP (wr. pulse time)	normal writes: 60 cycles (150 ns); 1.5x slow writes: 90 cycles (225 ns); 2.0x slow writes: 120 cycles (300 ns); 3.0x slow writes (default): 180 cycles (450 ns).
tCAS	1 cycle (2.5 ns)
endurance	normal writes: $5.000 * 10^6$ writes; 1.5x slow writes: $1.125 * 10^7$ writes; 2.0x slow writes: $2.000 * 10^7$ writes; 3.0x slow writes (default): $4.500 * 10^7$ writes.

and claim to be $1000\times$ faster and have $1000\times$ more endurance than NAND. Given the fact that commercial NAND Flashes provide 10^3 scale endurance (in P/E cycles) [33] and 100us scale write latency (program latency) [34], both Xpoint Memory and SCM may have an endurance of 10^6 scale and write latency of 100ns scale, which is in line with our experiment parameters.

We simulate multiple different memory write policies, as shown in Table 3. There are several configurations that can be mixed and matched— normal vs slow writes, eager write-backs, bank-aware mellow write-backs, write cancellation involving normal and/or slow writes, and wear quota. *BE-Mellow+SC+NC* means a system with both *Bank-Aware* and *Eager Mellow Writes*, and both normal writes and slow writes are cancellable; *BE-Mellow+SC+WQ* is the same except that only slow writes are cancellable and *Wear Quota* scheme is used. If without specific mentioning, we use default slow writes with $3.0\times$ latency in all the write policies.

Table 3: Memory Write Policies

Basic Policies	
Norm	Just using normal writes
Slow	Just using slow writes
B-Mellow	Using <i>Bank-Aware Mellow Writes</i>
BE-Mellow	Using both <i>Bank-Aware</i> and <i>Eager Mellow Writes</i>
E-Norm	Just using normal writes, but with eager writes
E-Slow	Just using slow writes, but with eager writes
Additional Write Choices	
+NC	Normal writes are cancellable
+SC	Slow writes are cancellable
+WQ	With <i>Wear Quota</i> scheme

Table 4: Single-Program Workloads and Their MPKI (Miss per 1000 Instructions) with a 2MB LLC

Workload	MPKI	Workload	MPKI	Workload	MPKI
leslie3d	5.95	hmmr	1.34	milc	19.49
GemsFDTD	15.34	zeusmp	4.53	mcf	56.34
libquantum	30.12	bwaves	5.58	lbm	31.72
stream	12.28	gups	8.91		

Table 5: Multi-Program Workloads

MIX1	hmmr zeusmp GemsFDTD libquantum
MIX2	libquantum GemsFDTD mcf bwaves
MIX3	mcf zeusmp hmmr leslie3d
MIX4	zeusmp milc bwaves lbm

In order to reduce the number of results displayed, we have chosen configurations that result in the best performance for the general option. For example, we show Norm and E-Norm+NC, but neither Norm+NC, nor E-Norm. This is because normal writes do not benefit enough from write cancellation to justify the drop in endurance. Write cancellation is important for eager write performance because it can avoid eager writes blocking the incoming reads. Also, the eager write queue does not trigger write drains, so cancelling eager slow writes will not increase the possibility of write drains.

For single-program workloads, we use 9 memory-intensive benchmarks from SPEC2006. To test the performance of our schemes under random and stream memory access patterns, we also include GUPS and `stream` benchmarks. We list these benchmarks in Table 4 with their MPKI (miss per 1000 Instructions). We warm up the cache for 6 billion instructions and then simulate in detail for another 2 billion instructions.

For multi-program workloads, we randomly pick 4 benchmarks from Table 4 and execute them together (Table 5). We warm up the cache until one program executes 6 billion instructions and then simulate in detail until one program executes 1 billion instructions.

6. RESULTS

In order to evaluate our *Mellow Writes* schemes, we first present the main tradeoff—performance vs lifetime. In order to better understand the reasons for these results, we break the performance down into bank utilization, write drain time, memory requests from the LLC, and memory requests issued to the memory banks. We also report the energy consumptions of the main memory system. We then provide a sensitivity study of several key parameters. Finally, we compare our best *Mellow Writes* with various kinds of static mechanisms. When not specifically stated, we use 3x write latency for all the slow writes.

6.1 Performance vs. Lifetime

We begin by presenting the fundamental tradeoff in our system: performance vs. lifetime. Figures 10 and 11 show the performance and lifetime, respectively, of our applications. We make the following observations:

- *E-Norm+NC* is designed for the highest performance. However, although it indeed gets the best performance for most of the benchmarks, it performs considerably worse than default (*Norm*) for 1bm (11% lower IPC) and also performs worse than *BE-Mellow+SC+WQ* for 1libquantum. This is because write cancellation may increase the write drain possibility by letting writes stay in write queue longer. Also *E-Norm+NC* has an unacceptable short lifetime. The lifetime suffers because of more write requests caused by eager write-backs and write cancellation. Therefore, eager writes and write cancellation should not be adopted with normal writes.
- *E-Slow+SC* has the longest lifetime. Unfortunately, the latency is far too high (geometric mean: 0.77x performance), with the worst being 0.46x (*lbm*). Even with write cancellation, implementing a system with only slow writes is not feasible.
- Using *Bank-Aware Mellow Writes (B-Mellow+SC)* improves the lifetime with negligible loss in performance.
- Combining *Bank-Aware Mellow Writes* with *Eager Mellow Writes (BE-Mellow+SC)* results in better performance and better lifetime than just using *Bank-Aware Mellow Writes (B-Mellow+SC)*. Overall, *BE-Mellow+SC* is a nice balance between lifetime and performance (Lifetime: 9.30 years; IPC: 1.06x of *Norm*).
- Although *Mellow Writes* schemes greatly improve the lifetime on average, some benchmarks achieve below the shortest acceptable lifetime (e.g., 8 years). In these cases, the *Wear Quota* scheme raises the lifetime to at least 8 years (see +*WQ* items in Figure 11). Not surprisingly, *Wear Quota* comes with performance cost. However, if we look at the three configurations with the *Wear Quota*, *BE-Mellow+SC+WQ* achieves the best performance. This is because, if the wear quota is reached, all schemes will have the same number of slow and normal writes, but the *Mellow Write* schemes will make better decisions about *which* writes should be slow vs normal. Therefore, for *Mellow Write* schemes, fewer sample periods are needed to use all slow writes than the normal scheme.

6.2 Bank Utilization

The utilization of a bank refers to the percentage of the time when corresponding bank is busy. We can see in Figure 12 that, not surprisingly, all configurations utilizing slow writes result in higher bank utilization. The mellow writes techniques (*B-Mellow+SC*, *E-Mellow+SC* and *BE-Mellow+SC*) sometimes (e.g., 1bm) result in higher utilization than globally slow writes with eager writes (*E-Slow+SC*), which may appear counter-intuitive. The reason for this phenomenon is that *E-Slow+SC* has considerably worse performance (e.g.,

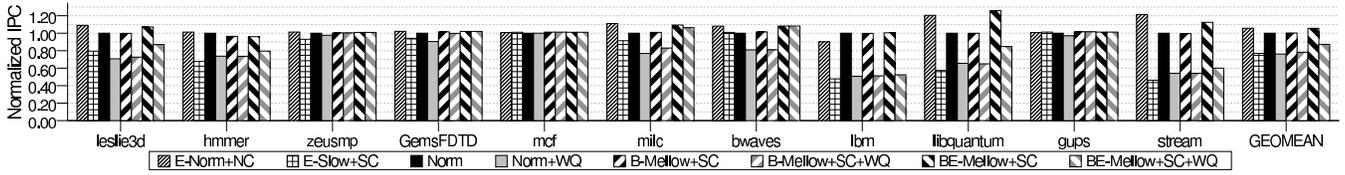


Figure 10: IPC (instruction per cycle) of systems with different write policies.

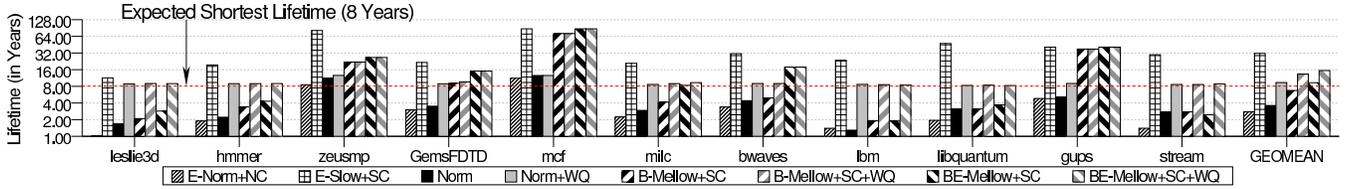


Figure 11: Resistive Memory Lifetime (in years, log scale) of systems with different write policies.

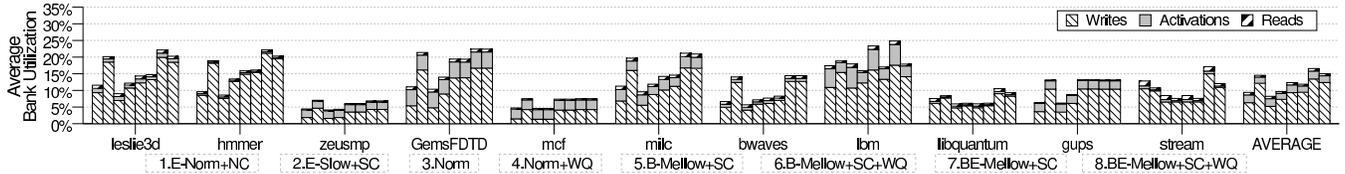


Figure 12: Average bank utilization of systems with different write policies.

Table 6: ReRAM cell parameters

	Read	Norm Set	Slow Set	Norm Reset	Slow Reset	
Voltage (V)	0.20	1.00	0.95	1.00	0.95	
Latency (nv)	–	150	450	150	450	
Power (uW)	0.02	–	–	–	–	
Energy per cell (pJ)	CellA	–	0.1	0.23	0.1	0.23
	CellB	–	0.2	0.46	0.2	0.46
	CellC	–	0.4	0.92	0.4	0.92
	CellD	–	0.8	1.84	0.8	1.84
	CellE	–	1.6	3.68	1.6	3.68

Table 7: Energy per operation of memristive main memory. We assume half of the bits in a write operation are subjected to the Set operation, and the other ones are subjected to Reset.

	Buffer Read (pJ)	Norm Write (pJ)	Slow Write (pJ)	Slow-Norm Write Ratio	Energy
CellA	1503.0	248.8	314.5	1.26	
CellB	1503.0	300.0	432.3	1.44	
CellC	1503.0	402.4	667.8	1.66	
CellD	1503.0	607.2	1138.8	1.88	
CellE	1503.0	1016.8	2080.9	2.05	

1bm) than mellow writes techniques, therefore fewer requests are sent to the memory controller in the same time period.

6.3 Write Drain Time

Write drains refer to the situation when the write queue occupancy reaches a threshold (usually full). When this occurs, the system prioritizes writes over reads until the write queue is drained. This is an expensive memory operation that directly impacts performance by delaying time-critical reads. This is also the main drawback of using slow writes: higher queue occupancy leading to more write drains. We study in detail of the write drain time in Figure 13.

When globally using the slow writes, the write queue fills often even if the *Eager Writes* are used (*E-Slow+SC*). *Bank Aware Mellow Writes* does not increase write drains (com-

pared with normal), and *Eager Mellow Writes* (*BE-Mellow+SC*) limit the write drain time to within 6% of the total execution time by proactively writing back data when the system is not busy. It is not surprising that using the *Wear Quota* scheme will increase the possibility of write drains. However, the percentages of write drain time in configurations with *Wear Quota* (*Norm+WQ*, *B-Mellow+SC+WQ* and *BE-Mellow+SC+WQ*) are still smaller than the write drain percentage when globally using slow writes (*E-Slow+SC*).

6.4 Memory Requests from LLC

We can see from Figure 14 that using *Eager Writes* transforms, on average, nearly half of the writes from normal writes to eager writes. Although *Eager Writes* may increase the number of write memory requests because of the inaccurate prediction of the unused blocks, this phenomenon is not obvious in our experiment (up to 2.2% increase of writes in benchmark *hmmer*). This reflects the fact that our mechanism to identify useless dirty blocks (described in Section 4.2.1) is relatively accurate.

6.5 Memory Requests to Memory Banks

Figure 15 shows the effects of write cancellation (which results in a second write after the read completes) and eager writebacks (which are wasted if the cacheline is modified again before eviction) on the number of issued memory requests to banks. *BE-Mellow+SC* issued substantially more requests to memory banks than *Norm*. However, as shown in Figure 14, the additional write requests due to eager writebacks are relatively few. Therefore, it is write cancellation which is mainly responsible for the increase of the issued writes to bank.

6.6 Energy Consumption of Main Memory

A potential drawback of a slow write is that it may consume more energy than a normal write. Therefore we simulate in detail the energy consumption of our schemes.

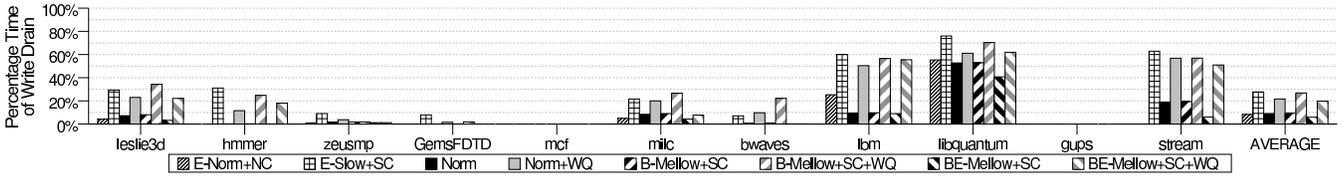


Figure 13: Percentage of time used by write drain operations.

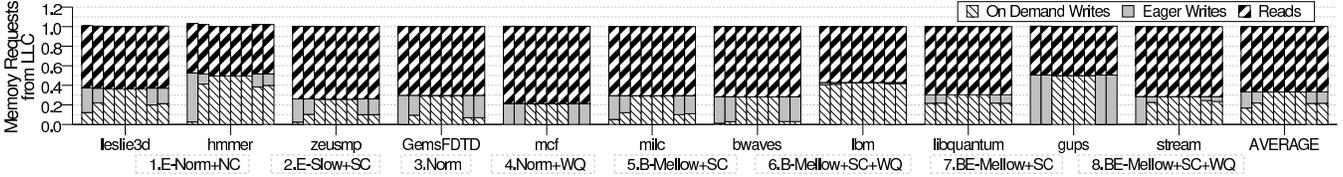


Figure 14: Number of memory requests from LLC to memory controller (normalized to the number of *Norm* policy).

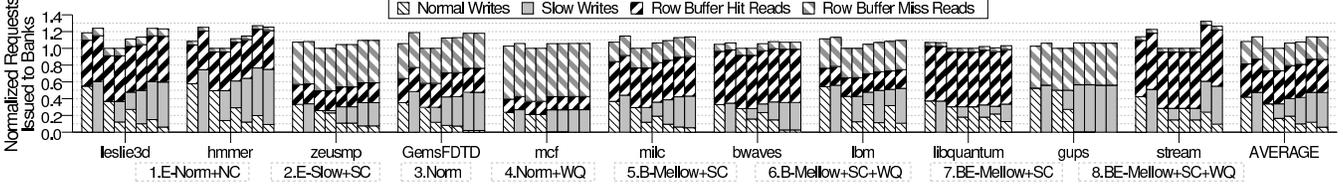


Figure 15: Number of memory requests issued to memory banks (normalized to the number of *Norm* policy).

The detailed energy simulation parameters are shown in Table 6. We use ReRAM under 22nm process. We assume that a $3\times$ slow write comes with $0.767\times$ dissipated power of a normal write, due to exponential dependence of ionic velocity on temperature [15]. Therefore, a slow write to the same ReRAM cell consumes $2.3\times$ energy of a normal write. Since the set/reset energy of a cell is a crucial design point of ReRAM, we model five different cells with their normal set/reset energy from 0.1pJ (*CellA*) to 1.6pJ (*CellE*). Then we use nvsim [35] to obtain the energy consumption of buffer read (in row buffer granularity) and normal/slow write (in cacheline granularity) operations of resistive main memory. We can see from Table 7 that, as the cell write energy decreases, the energy consumption difference of normal and slow writes also decreases—for *CellE* (1.6pJ/cell for normal set/reset), a slow write takes $2.05\times$ energy of a normal write; while for *CellA* (0.1pJ/cell for normal set/reset), a slow write only takes $1.26\times$ energy of a normal write. This is because the peripheral circuit also consumes energy while writing, when cell write energy decreases the energy consumption of peripheral circuit becomes more dominant.

We then calculate the whole main memory energy using the numbers from *CellC* in Table 7 and assume energy of a row-buffer hit read is 100pJ. The results are shown in Figure 16. On average, our best configuration (*BE-Mellow+SC+WQ*) consumes around $0.39\times$ more energy in main memory system than the default configuration (*Norm*). Given the fact that the main memory system usually consumes a relatively small portion of the whole system energy, the energy increase of *Mellow Writes* schemes is moderate/trivial compared to the total system energy consumption.

6.7 Multi-Program Workloads

Figure 17 shows the effectiveness of *Mellow Writes* for multi-program workloads. We can see that our schemes also work well in multi-program situations—compared to *Norm*, *BE-Mellow+SC* achieves $1.09\times$ IPC and $2.15\times$ lifetime.

6.8 Sensitivity to the Analytic Model

In Section 2, we show that if the write latency is slowed down by N times, an N^{Expo_Factor} times endurance can be achieved. And for resistive memory, an *Expo_Factor* of 2 is reasonable for our exploration. However, for various kind of non-volatile memory technologies, their *Expo_Factors* might be different. We investigate how our system performs as the exponential relationship between latency and lifetime changes. In addition to our default 2.0, we compare the lifetime of our schemes with four other *Expo_Factor* values, which are 1.0, 1.5, 2.5 and 3.0. We want to find out how dependent our results are on this exponential factor.

Figure 18 shows the results. Not surprisingly, the lifetime of both *Slow+SC* and *BE-Mellow+SC* improves when *Expo_Factor* goes up. However, the lifetime increase of *BE-Mellow+SC* is not as dramatic as for *Slow+SC*. The former gets around $0.5\times$ more lifetime for an *Expo_Factor* of 3.0 instead of 2.0, whereas the latter gets around $2\times$ more. This is because *Mellow Writes* schemes issue some amount of normal writes, and these normal writes contribute to a fixed amount of wear no matter how large *Expo_Factor* is. An important discovery is that, even for an *Expo_Factor* as pessimistic as 1.0, *BE-Mellow+SC* still gets a lifetime which is $1.47\times$ of the lifetime in a baseline system (*Norm*). Therefore, *Mellow Writes* are useful for a range of exponents, and thus technologies, and do not depend on a superlinear relationship between latency and endurance benefit.

6.9 Sensitivity to Bank-Level Parallelism

Here we demonstrate how available bank-level parallelism affects performance of *Mellow Writes*. Figure 19 shows the behavior of benchmark GemsFDTD with different numbers of banks. We can see from Figure 19 (a) that, as the number of banks decreases, the lifetime difference between *Norm* and *BE-Mellow+SC* diminishes, indicating the effectiveness of *Mellow Writes* diminishes. It is not surprising that this affects Bank-Aware mechanisms, since they depend on asym-



Figure 16: Energy consumption of main memory (see *CellC* in Table 7, normalized to *Norm*). Expected Shortest Lifetime (8 Years)

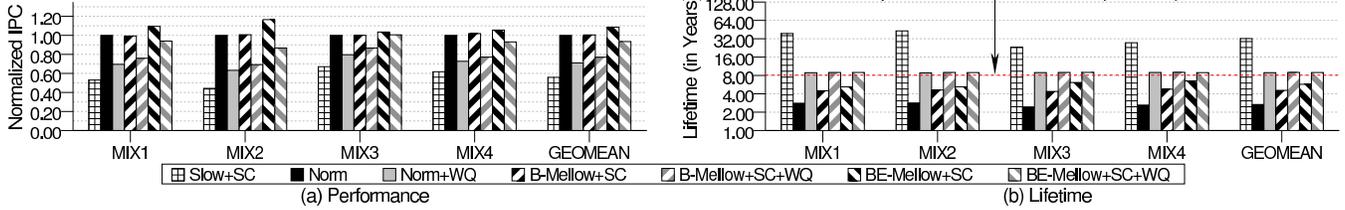


Figure 17: Effectiveness of *Mellow Writes* with multi-program workloads.

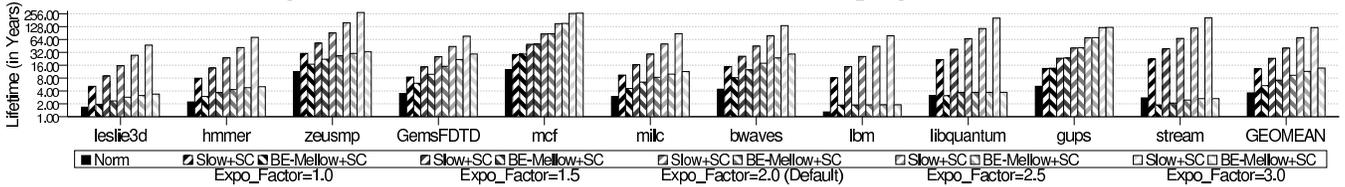


Figure 18: Sensitivity to exponent of latency to lifetime improvement relationship Equation 2.

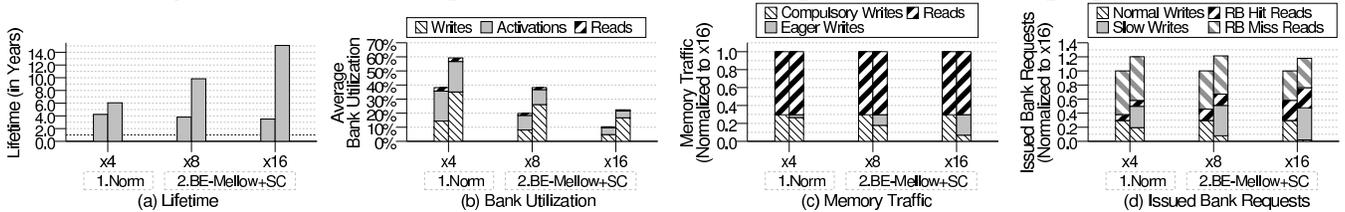


Figure 19: Sensitivity to bank-level parallelism of benchmark *GemsFDTD*.

metric use of banks. Figure 19(b) reveals the reason— as bank-level parallelism decreases, the utilization of each bank increases, leaving fewer intervals for eager and slow writes. Figure 19(c) clearly shows that the number of eager writes decreases dramatically as the number of banks decreases. Also, Figure 19(d) shows that the number of issued normal writes to the banks increases substantially as the number of banks decreases. Note that, although a substantial number of slow writes are issued with 4 banks, most of them get cancelled before finishing due to incoming reads.

Therefore, to ensure the effectiveness of *Mellow Writes*, it is important to run on a system with a sufficient number of banks.

6.10 Mellow Writes vs Static Policies

In Section 3, we show that it is hard to have a static mechanism (i.e., with fixed write latency and write policy) to fit different applications. To show the effectiveness of our *Mellow Writes Policy*. We compare *BE-MELLOW+SC+WQ* (our best *Mellow Writes* policy which can guarantee a minimal lifetime) with various kinds of static mechanisms. Since our modified *Eager Writes* policy can also be applied to a system with static write latency, we also include these static mechanisms (*E-Norm+NC* and *E-Slow+SC*) in our evaluation. We shown the results in Figure 20. For each benchmark, the column marked with red diamond is the best static policy (the

one guaranteeing minimal lifetime and delivering the best performance). We can see that, the best static mechanism is different for different benchmark, thus there is no single static mechanism which fits all applications.

As is shown in Figure 20, *BE-MELLOW+SC+WQ* successfully guarantees the minimal lifetime (in our experiment, 8 years) in all the applications. In 8 out of 11 applications, *BE-MELLOW+SC+WQ* outperforms or at equals the best static mechanism. We also investigate why our mechanism delivers worse performance in the rest of the benchmarks (i.e., *hmmer*, *lbm* and *stream*)—it turns out these benchmarks are very sensitive to write latency in some cases. A possible modification for this situation is to adopt multiple write latencies instead of just two in our schemes. In this case, deciding which write latency to use is a major challenge, and this will be our future work.

Overall, *BE-MELLOW+SC+WQ* shows a nice balance between lifetime and performance requirements for all the benchmarks, and such a balance cannot be fulfilled with a single static mechanism.

7. RELATED WORK

Wear Leveling and Limiting. Many techniques exist for implementing wear-leveling for non-volatile memories. Start-Gap [23] employs a novel shift-based approach that is able to achieve 95% of ideal memory lifetime with only 8

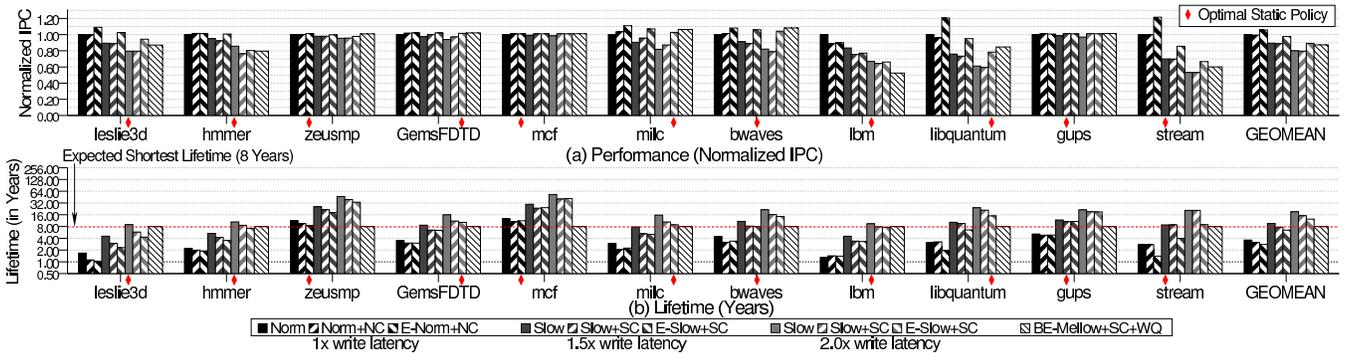


Figure 20: Compare *BE-Mellow+SC+WQ* with various kinds of static policies. For each benchmark, the column with red diamond is best static mechanism (i.e., the static mechanism which guarantees 8 years lifetime and delivers the best performance). We can see that there is no single static mechanism which suits all the benchmarks, and our *BE-Mellow+SC+WQ* outperforms or equals to the best static mechanism in 8 out 11 benchmarks.

bytes of storage overhead. We use Start-Gap in our system. Other shift-based approaches exist, such as shifting cache lines with a page [1], and shifting bits in a line, or lines in a segment [36]. Security Refresh [37] uses randomized address mappings within a bank for distribute wear as well as prevent malicious wear-out attacks. Well-known techniques exist for wear-limiting, such as DRAM buffering [1]. Others, such as Flip-N-Write [2], exploit specific properties of the data being written to reduce the number of writes on a per-cell basis. Similar to DRAM buffering, Saadeldeen *et al.* [38] also use a small SRAM buffer in their ReRAM-based branch prediction scheme. All such techniques can be classified as *physical* techniques, in that they alter the actual contents of the memory in order to reduce wear. Our mellow writes concept is orthogonal to these, as it uses *temporal* properties of write operations to reduce wear.

Latency-Density Tradeoffs in MLC NVMs. It is well known that, when using MLC (multi-level cell) NVM, there is a trade-off between write/read latency and storage density. Prior proposals try to balance such a trade-off in main memory [39], file storage [40] and coherence directories [41] by adaptively using fast and slow operations. The basic notion of these proposals is to use fast(single-level) accesses for performance critical read/writes operations, and slow(multi-level) accesses for other ones. In this work, we also follow this intuition to avoid performance loss by not using slow writes in performance-critical situations.

Write Cancellation and Eager Writeback. To avoid of performance loss due to long latency write in NVM, Qureshi *et al.* [18] introduce the concept of *write cancellation* (also known as *read preemption* [19]) which services the incoming reads immediately by cancelling conflict writes. Lee *et al.* [42] propose a method by which LRU dirty cache lines are written back before eviction, called *eager writeback*. Doing so reduces memory bandwidth contention between demand reads and writebacks, thus improving system performance. Qureshi *et al.* [43] also propose a scheme to improve PCM performance by early and eagerly writing back the long latency SET operations, which can be viewed as a variation of *eager writeback*. Both write cancellation and a modified version of eager writeback are integral parts of our mellow writes concept. Eager writeback provides the mechanism by

which some slow writes are generated, and write cancellation prevents these slow writes from impacting system performance.

Cache Management. Qureshi *et al.* [21] propose a method of partitioning an LRU-policy cache among concurrently executing applications. This work provides the motivation for our Eager Mellow Writes writeback criteria, in that we identify writeback candidates based on their utility rather than simply their LRU stack position.

8. CONCLUSIONS

In this paper we explore the trade-off that, for non-volatile memories, there exists a linear to cubic (representatively to be quadratic for resistive memories) endurance advantage to performing writes slowly using a smaller write dissipated power. Although slower writes can dramatically improve lifetime, they may also degrade overall system performance.

To utilize the lifetime benefit of slow writes as well as avoid their performance penalty, we introduce two schemes that issue slow writes when the memory system is not busy, namely *Bank-Aware Mellow Writes* and *Eager Mellow Writes*. The *Bank-Aware Mellow Writes* scheme performs slow writes only when there is no other request to the same bank. The *Eager Mellow Writes* scheme eagerly and slowly writes back the dirty blocks in the LLC which are predicted to have a low probability of being re-written. In addition, we also introduce *Wear Quota* scheme to ensure the minimal expected memory lifetime—when the *Wear quota* of bank in all previous sample periods is reach, the bank can only issue slow writes in the coming period. Our experiments show that a combination of these two *Mellow Writes* schemes extends lifetime with minimal performance impact.

This combination can achieve $2.58\times$ in lifetime and $1.06\times$ in performance of a baseline system using normal write speed. Meanwhile, *Wear Quota* is a safe scheme which can guarantee the minimal expected lifetime (e.g., 8 years) with relatively small performance loss. Finally, we find that while the lifetime benefit of Mellow Writes schemes increases as the endurance benefit of slower writes increases, they are advantageous even with a linear relationship. This makes our design and analysis applicable to a wide range of non-volatile memory technologies.

9. REFERENCES

- [1] M. K. Qureshi, V. Srinivasan, and J. A. Rivers, "Scalable high performance main memory system using phase-change memory technology," in *Proceedings of the 36th Annual International Symposium on Computer Architecture*, pp. 24–33, 2009.
- [2] S. Cho and H. Lee, "Flip-n-write: A simple deterministic technique to improve pram write performance, energy and endurance," in *Microarchitecture, 2009. MICRO-42. 42nd Annual IEEE/ACM International Symposium on*, pp. 347–357, 2009.
- [3] D. B. Strukov, "Endurance–write speed tradeoffs in nonvolatile memories." To appear in *Applied Physics A*, and a preprint is available in: <http://arxiv.org/abs/1511.07109>.
- [4] X. Liu, V. Patel, Z. Tan, K. K. Likharev, and J. E. Lukens, "High-quality aluminum-oxide tunnel barriers for scalable, floating-gate random-access memories (fgam)," in *Proc. Int. Conf. on Memory Technology and Design (ICMTD)*, pp. 235–237, 2007.
- [5] H.-C. Yu, K.-C. Lin, K.-F. Lin, C.-Y. Huang, Y.-D. Chih, T.-C. Ong, J. Chang, S. Natarajan, and L. Tran, "Cycling endurance optimization scheme for 1mb stt-mram in 40nm technology," in *Solid-State Circuits Conference Digest of Technical Papers (ISSCC), 2013 IEEE International*, pp. 224–225, 2013.
- [6] K. K. Likharev, "Layered tunnel barriers for nonvolatile memory devices," *Applied Physics Letters*, vol. 73, no. 15, pp. 2137–2139, 1998.
- [7] M. D. Pickett, D. B. Strukov, J. L. Borghetti, J. J. Yang, G. S. Snider, D. R. Stewart, and R. S. Williams, "Switching dynamics in titanium dioxide memristive devices," *Journal of Applied Physics*, vol. 106, no. 7, p. 074508, 2009.
- [8] J. McPherson, J.-Y. Kim, A. Shanware, and H. Mogul, "Thermochemical description of dielectric breakdown in high dielectric constant materials," *Applied Physics Letters*, vol. 82, no. 13, pp. 2121–2123, 2003.
- [9] J. J. Yang, D. B. Strukov, and D. R. Stewart, "Memristive devices for computing," *Nature Nanotechnology*, vol. 8, no. 1, pp. 13–24, 2013.
- [10] K. Likharev, "Electronics below 10 nm," in *In Nano and Giga Challenges in Microelectronics*, pp. 27–68, 2003.
- [11] V. G. Karpov, Y. A. Kryukov, I. V. Karpov, and M. Mitra, "Field-induced nucleation in phase change memory," *Phys. Rev. B*, vol. 78, no. 5, p. 052201, 2008.
- [12] S. Raoux, D. Ielmini, M. Wuttig, and I. Karpov, "Phase change materials," *MRS bulletin*, vol. 37, no. 02, pp. 118–123, 2012.
- [13] E. Tsymal, A. Gruverman, V. Garcia, M. Bibes, and A. Barthélémy, "Ferroelectric and multiferroic tunnel junctions," *MRS bulletin*, vol. 37, no. 2, pp. 138–143, 2012.
- [14] R. Waser, R. Dittmann, G. Staikov, and K. Szot, "Redox-based resistive switching memories–nanoionic mechanisms, prospects, and challenges," *Advanced Materials*, vol. 21, no. 25–26, pp. 2632–2663, 2009.
- [15] D. B. Strukov and R. S. Williams, "Exponential ionic drift: fast switching and low volatility of thin-film memristors," *Applied Physics A*, vol. 94, no. 3, pp. 515–519, 2009.
- [16] N. Mott and R. Gurney, *Electronic Processes in Ionic Crystals*. Dover, New York, 2nd edn ed., 1940.
- [17] N. F. Mott and R. W. Gurney, "Electronic processes in ionic crystals," 1948.
- [18] M. Qureshi, M. Franceschini, and L. Lastras-Montano, "Improving read performance of phase change memories via write cancellation and write pausing," in *High Performance Computer Architecture (HPCA), 2010 IEEE 16th International Symposium on*, pp. 1–11, 2010.
- [19] G. Sun, X. Dong, Y. Xie, J. Li, and Y. Chen, "A novel architecture of the 3d stacked mram l2 cache for cmps," in *High Performance Computer Architecture, 2009. HPCA 2009. IEEE 15th International Symposium on*, pp. 239–249, 2009.
- [20] R. L. Mattson, J. Gecsei, D. R. Slutz, and I. L. Traiger, "Evaluation techniques for storage hierarchies," *IBM Systems journal*, vol. 9, no. 2, pp. 78–117, 1970.
- [21] M. K. Qureshi and Y. N. Patt, "Utility-based cache partitioning: A low-overhead, high-performance, runtime mechanism to partition shared caches," in *Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 423–432, 2006.
- [22] M. K. Qureshi, A. Jaleel, Y. N. Patt, S. C. Steely, and J. Emer, "Adaptive insertion policies for high performance caching," in *Proceedings of the 34th Annual International Symposium on Computer Architecture*, pp. 381–391, 2007.
- [23] M. K. Qureshi, J. Karidis, M. Franceschini, V. Srinivasan, L. Lastras, and B. Abali, "Enhancing lifetime and security of pcm-based main memory with start-gap wear leveling," in *Proceedings of the 42Nd Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 14–23, 2009.
- [24] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. D. Hill, and D. A. Wood, "The gem5 simulator," *SIGARCH Comput. Archit. News*, vol. 39, no. 2, pp. 1–7, 2011.
- [25] M. Poremba and Y. Xie, "Nvmain: An architectural-level main memory simulator for emerging non-volatile memories," in *VLSI (ISVLSI), 2012 IEEE Computer Society Annual Symposium on*, pp. 392–397, 2012.
- [26] C. Xu, D. Niu, N. Muralimanohar, R. Balasubramonian, T. Zhang, S. Yu, and Y. Xie, "Overcoming the challenges of crossbar resistive memory architectures," in *High Performance Computer Architecture (HPCA), 2015 IEEE 21st International Symposium on*, pp. 476–488, Feb 2015.
- [27] J. Choi, J.-S. Kim, I. Hwang, S. Hong, S. Jeon, S.-O. Kang, B. Park, D. Kim, M. Lee, and S. Seo, "Different resistance switching behaviors of nio thin films deposited on pt and srro3 electrodes," *Applied Physics Letters*, vol. 95, no. 2, p. 022109, 2009.
- [28] I. H. Inoue, S. Yasuda, H. Akinaga, and H. Takagi, "Nonpolar resistance switching of metal/binary-transition-metal oxides/metal sandwiches: Homogeneous/inhomogeneous transition of current distribution," *Physical Review B*, vol. 77, no. 3, p. 035105, 2008.
- [29] R. Waser, R. Dittmann, G. Staikov, and K. Szot, "Redox-based resistive switching memories–nanoionic mechanisms, prospects, and challenges," *Advanced Materials*, no. 21, pp. 2632–2663, 2009.
- [30] M.-J. Lee, C. B. Lee, D. Lee, S. R. Lee, M. Chang, J. H. Hur, Y.-B. Kim, C.-J. Kim, D. H. Seo, S. Seo, et al., "A fast, high-endurance and scalable non-volatile memory device made from asymmetric ta2o5-x/tao2- x bilayer structures," *Nature materials*, vol. 10, no. 8, pp. 625–630, 2011.
- [31] "Intel and Micron produce breakthrough memory technology." http://newsroom.intel.com/community/intel_newsroom/blog/2015/07/28/intel-and-micron-produce-breakthrough-memory-technology.
- [32] "Sandisk and HP launch partnership to create memory-driven computing solutions." <https://www.sandisk.com/about/media-center/press-releases/2015/sandisk-and-hp-launch-partnership>.
- [33] E. Grochowski and R. E. Fontana Jr, "Future technology challenges for nand flash and hdd products," *Flash Memory Summit*, 2012.
- [34] L. M. Grupp, A. M. Caulfield, J. Coburn, S. Swanson, E. Yaakobi, P. H. Siegel, and J. K. Wolf, "Characterizing flash memory: anomalies, observations, and applications," in *Microarchitecture, 2009. MICRO-42. 42nd Annual IEEE/ACM International Symposium on*, pp. 24–33, IEEE, 2009.
- [35] X. Dong, C. Xu, Y. Xie, and N. P. Jouppi, "Nvsm: A circuit-level performance, energy, and area model for emerging nonvolatile memory," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, vol. 31, no. 7, pp. 994–1007, 2012.
- [36] P. Zhou, B. Zhao, J. Yang, and Y. Zhang, "A durable and energy efficient main memory using phase change memory technology," in *Proceedings of the 36th Annual International Symposium on Computer Architecture*, pp. 14–23, 2009.
- [37] N. H. Seong, D. H. Woo, and H.-H. S. Lee, "Security refresh: Prevent malicious wear-out and increase durability for phase-change memory with dynamically randomized address mapping," in *Proceedings of the 37th Annual International Symposium on Computer Architecture*, pp. 383–394, 2010.
- [38] H. Saadelddeen, D. Franklin, G. Long, C. Hill, A. Browne, D. Strukov, T. Sherwood, and F. T. Chong, "Memristors for neural branch prediction: a case study in strict latency and write endurance challenges," in *Proceedings of the ACM International Conference on Computing Frontiers*, pp. 26:1–26:10, 2013.

- [39] M. K. Qureshi, M. M. Franceschini, L. A. Lastras-Montaña, and J. P. Karidis, "Morphable memory system: A robust architecture for exploiting multi-level phase change memories," in *Proceedings of the 37th Annual International Symposium on Computer Architecture*, pp. 153–162, 2010.
- [40] X. Dong and Y. Xie, "Adams: Adaptive mlc/slc phase-change memory design for file storage," in *Design Automation Conference (ASP-DAC), 2011 16th Asia and South Pacific*, pp. 31–36, 2011.
- [41] L. Zhang, D. Strukov, H. Saadeldeen, D. Fan, M. Zhang, and D. Franklin, "Spongedirectory: Flexible sparse directories utilizing multi-level memristors," in *Proceedings of the 23rd International Conference on Parallel Architectures and Compilation*, pp. 61–74, 2014.
- [42] H.-H. S. Lee, G. S. Tyson, and M. K. Farrens, "Eager writeback - a technique for improving bandwidth utilization," in *Proceedings of the 33rd Annual ACM/IEEE International Symposium on Microarchitecture*, pp. 11–21, 2000.
- [43] M. K. Qureshi, M. M. Franceschini, A. Jagmohan, and L. A. Lastras, "Preset: Improving performance of phase change memories by exploiting asymmetry in write times," in *Proceedings of the 39th Annual International Symposium on Computer Architecture*, pp. 380–391, 2012.