



lower switching current than conventional MRAMs (100-200uA v.s. 10mA [2,12,21]), which makes it feasible for cache or memory. Furthermore, STT-RAM's threshold current decreases as the size of MTJ becomes smaller, resulting in better scalability than previous MRAM designs [2]. However, STT-RAM still faces challenge of high write energy, compared to an SRAM cell write [2, 14].

A typical design of STT-RAM cell array is shown in Figure 3 [2, 12]. An NMOS connected to word line (WL) is used to select a row of cells, and voltage is applied between bit line (BL) and source line (SL) to perform read/write operations:

- When reading from the cell, a small negative voltage (usually -0.1V) is applied between BL and SL [2, 12]. The amount of current flowing through the cell is dependent on the resistance of MTJ, which is sensed by a sense amplifier to output the stored data.
- When writing a '0' to the cell, a positive voltage (typically 1.2V) is applied between SL and BL, creating a current flow from SL to BL. When writing a '1', a negative voltage (typically -1.0V) is applied between SL and BL, creating a current flow in the opposite direction. The write current is significantly larger than the read current and the duration is also much longer [2, 3].

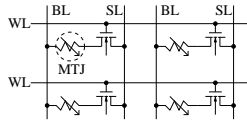


Figure 3: Typical design of STT-RAM cell array [2, 12].

### 3. EARLY WRITE TERMINATION (EWT)

#### 3.1 The Opportunity

Several existing studies have shown that there is a high probability that a write into a cache or memory location does not change its content, and therefore can be removed. Such an observation has been used at the word level for L1 cache [8], multiprocessors [9], and off-chip memories [7, 20]. Our evaluation demonstrates that this phenomenon is more significant at the bit level, and can be used to throttle those bit-writes that do not change the stored value. Fig. 4 shows the results of our evaluation for a 16MB STT-RAM L2 cache (experimental settings will be discussed in Section 5.1). On average, we see that about 88% of bit-writes are redundant, which implies a significant amount of removable bit writes and a great potential of energy savings in a STT-RAM cache.

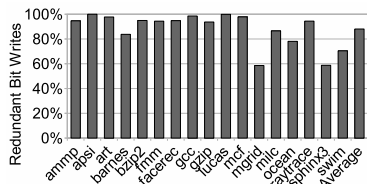


Figure 4: Redundant bit writes in 16MB STT-RAM L2 cache.

#### 3.2 The Rationale

To remove the unnecessary writes, one needs to know if the new bit value is different from the old bit value. Since our STT-RAM is at the L2 level, one simple way is to utilize the "dirty" bit in L1 to pass or terminate the writes into the L2, analogous to the write removal technique at the memory level [7]. That is, each L1 word is tagged with a dirty bit, and only dirty words are written back to L2 upon eviction. However, this increases the number of dirty bit management in L1 (one bit per word vs. one bit per line by default). More importantly, it cannot capture all the opportunities in write removal. For example, an L1 word being dirty does not warrant that the value is different from what is already stored in L2. Also, even if the value is different, there still might be many bits that are the same. Hence, this technique cannot bring significant reduction in write energy.

To exploit the full opportunity in redundant bit-writes, we could first read out the cache content, compare it with the new values, and write back only the different bits. However, this method entails that every write is preceded with a read operation. Although reads consume much less energy and are much faster than writes in STT-RAM, we could do better than the above method based on the following unique features of STT-RAM cells:

1. When writing to an STT-RAM cell, the change of MTJ resistance is not a gradual procedure. Instead, resistance changes abruptly near the end of a write cycle [1]. This means that at the early stage of a write operation, STT-RAM cell still holds its valid old value.
2. A read operation is performed through applying a voltage between the bit-line and source line, followed by sensing the resulting current on the two lines to determine the MTJ's resistance. This can be performed during the write operation as the latter simply keeps the voltage for a sustained amount of time to change the MTJ state.
3. A write operation in STT-RAM is much longer than a read. A typical write pulse of a STT-RAM cell is 10ns, below which the switching current increases rapidly [2,3]. However, reading from a cell can complete in less than 1ns even in a large L2 [2]. This not only gives us adequate room to sense the old value during the early stage of a write, but also implies great opportunity in saving energy by terminating the write current as soon as redundancy is detected.

Based on those observations, we propose a novel write scheme with the capability of early termination in case of a redundant write. The main goal is to greatly reduce the write energy without impact on performance. The basic idea is to sample the resistance of the MTJ (old value) at early stage of a write operation, and throttle the write current if old value is the same as the new value. To achieve this goal, we need an additional sense amplifier  $SA_w$  for each column to sense the resistance of the MTJ during a write operation, and some additional logic to generate the control signals. The reason we use a separate sense amplifier is because the original sense amplifier is used only for read operations which work with only small currents. Hence, it cannot be used for write operations which generate much larger current than reads. As we will show later, we use a normal voltage sense amplifier for  $SA_w$  due to its simplicity and low power. The procedure of EWT is as follows.

- When a write operation begins, write voltage is applied between BL and SL to form a write current.
- When the signals are stabilized, sense amplifier  $SA_w$  is enabled to sense the resistance of the MTJ (i.e. old value).
- After the old value is sensed, it is saved in a latch and  $SA_w$  is turned off to save power.
- If the old value is the same as the new value, a control signal  $WCUT$  is generated to shut off the write circuit and terminate the operation on this cell. Otherwise, write operation on this cell continues normally.

As we can see, the above process does not require an extra read to precede a write because sensing the old value is done together with the write operation. Comparing to a previous scheme used in Phase Change Memory that mandates a read before a write [20], EWT does not introduce any overhead in performance. In fact, our experiments show that EWT sometimes even improves performance a little because some write requests (an entire L1 cache line) can be completely throttled and terminated in their early stage.

#### 3.3 The Circuit Design

Our EWT implementation does not require any change to existing read or write circuit. Instead, it is designed to work as an "add-on" to a working MRAM read and write circuit. Therefore, EWT can be easily integrated into existing designs.

The write operation starts with applying a positive voltage, for writing a '0', or negative voltage, for writing a '1', between the SL and BL. We add a pass gate on both BL and SL, as shown in Figure 5. These pass gates serve two purposes: 1) when it is detected that the write is redundant, the pass gates are turned off to cut the write current on BL and SL; 2) they are also small loads added to

the BL and SL to show the write path voltage distribution that is determined by the MTJ's resistance. This is used to detect the stored value in MTJ.

For example, when a '0' is written, a positive voltage is applied between SL and BL creating a current flow from SL to BL. Hence, there is a voltage drop on the pass gate on the SL. However, the magnitude of this drop is determined mainly by the resistance of the MTJ (other wire loads are relatively constant). If it is storing a '1', meaning that the resistance is high, the voltage drop on the pass gate is relatively small and  $V_{in0}$  is relatively high. On the other hand, if the MTJ is storing a '0',  $V_{in0}$  is relatively low. Such a voltage difference is used to determine the stored value in MTJ.

To avoid disturbance on write current and ensure that  $V_{in0}$  can be correctly sensed, a conversion circuit is used to magnify  $V_{in0}$  into  $V_{sense}$ , as illustrated in Figure 6. The circuit is similar to a basic differential amplifier [22]. In this conversion circuit,  $P_1$ ,  $P_2$  and  $N_2$  form a current mirror to provide output current on  $N_1$ . They are properly sized so that output current is large enough for sensing. Signal  $SE$  is used to turn on and off the conversion circuit. Input voltage  $V_{in}$  is connected to  $N_1$ 's gate, which controls its equivalent resistance. For example, if  $V_{in}$  is relatively higher,  $N_1$ 's equivalent resistance is smaller, and vice versa. This difference is reflected at  $V_{sense}$  as output current flows through  $N_1$ . In this way, the difference on  $V_{in}$  is magnified into the difference on  $V_{sense}$ , which is used as the input to the sense amplifier referred in Figure 5.

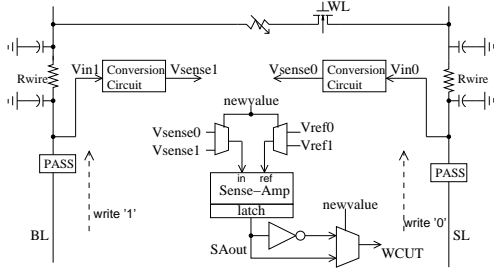


Figure 5: EWT circuit design in a column.

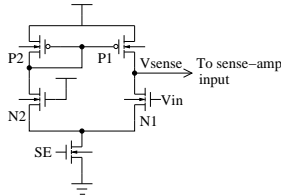


Figure 6: The conversion circuit.

Writing a '1' is carried by applying a negative voltage between SL and BL, and the process is opposite to writing a '0'. Therefore, we used a symmetric conversion circuit on the BL for writing a '1'. The two outputs,  $V_{sense0}$  and  $V_{sense1}$  are then sent to the sense amplifier  $SA_w$  through a mux controlled by the new bit value, e.g.  $V_{sense0}$  is selected if a '0' is written. At the same time, two reference voltages  $V_{ref0}$  and  $V_{ref1}$  are also sent to  $SA_w$  for comparing with  $V_{sense0}$  and  $V_{sense1}$  respectively. We need different reference voltages because the circuit is not entirely symmetric [1, 12] and variation of  $V_{sense}$  is also different on BL and SL.

The last step is to generate a control signal  $WCUT$  to either turn off the write circuit or let it continue, based on the result of the sense amplifier  $SA_{out}$ . The  $SA_{out}$  indicates whether the MTJ is storing a '0' or a '1'. We would need to compare it with the new value to derive  $WCUT$ . If they are equal,  $WCUT$  is high, otherwise,  $WCUT$  is low. Fortunately, we do not need to use a comparison circuit such as a XOR/XNOR gate to implement it. Table 1 summarizes the signals and actions during a write. As we can see that  $WCUT = SA_{out}$  when writing a '0', and  $WCUT = \overline{SA_{out}}$  when writing a '1'. Hence, either  $SA_{out}$  or  $\overline{SA_{out}}$  is used to turn off the pass gates and the write circuit. This can be implemented by an inverter and mux, which is smaller and simpler than XOR/XNOR gate. In our HSPICE simulation,  $WCUT$  can be

generated in 0.536ns in a 16MB L2 cache after the word-line is selected. Therefore, redundant bit writes can be detected and throttled at a very early stage.

Table 1: Sensing Signals

Old	New	$V_{in0}$	$V_{sense0}$	$SA_{out}$	Action
0	0	Lower	$> V_{ref0}$	1	Cut
1	0	Higher	$< V_{ref0}$	0	Continue
Old	New	$V_{in1}$	$V_{sense1}$	$SA_{out}$	Action
0	1	Lower	$> V_{ref1}$	1	Continue
1	1	Higher	$< V_{ref1}$	0	Cut

### 3.4 Overhead

We implemented the EWT circuit and simulated them in HSPICE using 45nm technology. We measured the additional energy introduced by EWT circuits, including pass gates, conversion circuit, muxes, sense amplifier, latch, and inverters. These components are added on per column basis. That is, all cells in one column share one set of the EWT circuit. We measured that the average energy overhead per cell per write is 89.29fJ. Comparing to the 2.767pJ cell write energy (discussed in Section 4.2), the energy overhead introduced by EWT is 3.23%.

The estimated area introduced by EWT circuit is about  $13.44\mu m^2$  per column (330 $\mu m$  long column). Comparing to the total area of a 16MB STT-RAM cache (calculated by CACTI), the estimated area overhead is 4.17%.

Since EWT is carried within a write operation, there is no performance overhead to the write latency. On the contrary, some write requests can even finish earlier if all the bits are the same as what are already stored in the cache. This leads to a slight performance gain, which we will see in Section 5.2.

## 4. MODELING STT-RAM AND EWT

To measure how much energy savings we can achieve through our EWT design, we first modeled a STT-RAM L2 cache in both performance and energy, and then compared it to a baseline STT-RAM cache without the EWT. As we have shown in Figure 3, STT-RAM cache uses similar peripheral logic as a regular SRAM cache. Hence, we use the existing cache modeling tool CACTI [16] to derive both the latency and the energy values for the peripheral logic such as the H-tree, decoder, word-line, bit-line, sense amplifiers etc., and combine them with the STT-RAM cell's latency and energy. We chose the 45nm technology library, and the low operation power (LOP) peripheral design due to the high dynamic power required by STT-RAM cells.

### 4.1 Latency

The breakdown of read and write component latencies are listed in Table 2. The latencies are rounded up to CPU cycles when used in our simulator. We refer to a recent work on STT-RAM cache [2] for STT-RAM's cell latency. For read operation, this is essentially the sense-amplifier delay, which is assumed to be 20% slower than SRAM's sense amplifier [2].

For write operations, we used a 10ns pulse width as mentioned earlier. However, if the entire write access (a cache line) is throttled, the write pulse width is equal to the time required for redundancy detection which is 0.536ns, as measured from our HSPICE simulation. Therefore, a write with EWT may take shorter time than in the baseline.

Table 2: Per-Access Read/Write Latency

	Read	Write (Base)	Write (EWT)
H-tree in	2.010ns	2.010ns	2.010ns
Word-line + Decoder	0.544ns	0.544ns	0.544ns
Bit-line	0.800ns	N/A	N/A
Sense-amp	1.006ns	N/A	N/A
H-tree out	1.872ns	N/A	N/A
Write Pulse	N/A	10ns	10ns / 0.536ns
Total	6.232ns	12.554ns	12.554ns / 3.090ns

### 4.2 Dynamic Energy

The breakdown of dynamic energy for reads and writes are shown in Table 3–4. For dynamic cell energies, we referred to the results from the recent work [2], and scale them to 45nm technology.

When EWT is enabled, the write energy is no longer a fixed value. Instead, it is the sum of three parts: peripheral energy  $E_{peripheral}$ , overhead  $E_{overhead}$  and a varying cell energy  $E_{cells}$  due to a value change:

$$E_{EWTwrite} = E_{peripheral} + E_{overhead} + E_{cells}$$

$E_{peripheral}$  is the energy consumed by the peripheral logic. This is 0.203nJ, same as in baseline.  $E_{overhead}$  is the energy consumed by the EWT circuits. This part is 0.0457nJ per write access, calculated as per cell overhead multiplied by the number of cells in a cache line (512 in our case) since there is one set of EWT circuit per column.  $E_{cells}$  is the energy required by those cells that are updated. This variable part depends on how many cells are actually changed in a write request. It can be expressed as:

$$E_{cells} = N_{changed} \times E_{cellchange} + N_{unchanged} \times E_{unchanged}$$

Where  $E_{cellchange}$  is the energy used to change one cell, which is 2.767pJ in our model. This is obtained from scaling the results in [2] to 45nm technology. Write operations on unchanged cells are terminated at the end of 0.536ns, which amounts to 0.148pJ per cell for  $E_{unchanged}$ . In summary, per-access write energy with EWT can be expressed as:

$$E_{EWTwrite} = E_{peripheral} + E_{overhead} + N_{changed} \times 2.767pJ + N_{unchanged} \times 0.148pJ$$

**Table 3: Per-Access Read/Write Energy**

	Read	Write (Base)	Write (EWT)
Peripheral	0.192nJ	0.203nJ	0.203nJ
Overhead	N/A	N/A	0.0457nJ
Cells	0.013nJ	1.417nJ	Variable (Table 4)
Total	0.205nJ	1.620nJ	Variable (Table 4)

**Table 4: EWT Write Energy (Per-Cell)**

$E_{cellchange}$	2.767pJ per cell
$E_{unchanged}$	0.148pJ per cell

### 4.3 Leakage Energy

Since STT-RAM has negligible leakage in the cell, peripheral leakage becomes the dominant part in the total leakage of the cache. We used CACTI to estimate the leakage power for the peripheral logic. This value is 0.265W for the 16MB STT-RAM cache used in our experiment. Due to the non-volatile nature of STT-RAMs, we can further reduce the leakage power by power gating the cache banks if they are idle. This significantly reduces the leakage power of a STT-RAM cache. CACTI results show that the leakage power of an idle bank in our model is 0.388mW, in a 16-bank 16MB cache. In addition, we assume that there is a 1ns delay to power on a bank. Such leakage saving scheme is applied to both the baseline and our design with EWT.

## 5. EVALUATIONS

### 5.1 Experimental Setup

We used a simulator based on Simics [11] to simulate a 3D architecture with a 16MB L2 STT-RAM stacked on top of a 4-core CMP, the most feasible way to integrate STT-RAMs with CMOS technology as studied previously [14]. Core frequencies are set to 1GHz due to thermal constraint in a 3D architecture. We enhanced the cache model in Simics to model the energy and delay of STT-RAM, as well as the contention on cache banks and wake-up delays of idle banks.

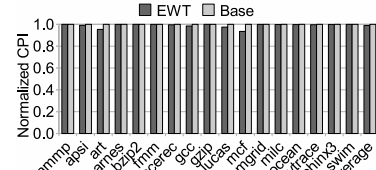
We use a variety of workloads from SPLASH2 [13], SPEC2K and SPEC2006. Both memory intensive and computational workloads are evaluated in our experiments. Results are collected through execution-driven simulations. For each workload, we skipped its initialization phase, warmed up for 50M instructions and ran for 100M instructions.

**Table 5: Simulation Parameters**

Processor core	4 cores, each core runs at 1GHz
L1 Cache	Private L1 cache (32K I-cache and 32K D-cache), 64-byte lines, 4-way set associative, 3 cycles access time
L2 Cache	Shared L2 cache (STT-RAM), 16MB, 64-byte lines, 16-way set associative, 16 banks
Main Memory	4GB memory, 50 cycles access time

### 5.2 Performance

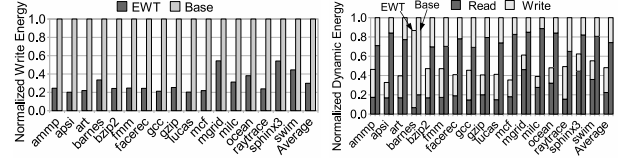
As we discussed previously, EWT does not introduce any performance penalty to cache accesses. Instead, write requests may finish early if no bit change is needed. Therefore, EWT can reduce average write latency and the contention on cache banks, which results in slight improvement in performance. Figure 7 shows our results in Cycles Per Instruction (CPI), normalized to the baseline. We observe a 3%-7% of CPI improvement in memory intensive workloads such as *mcf*, *art*, *lucas*, and the average CPI improvement over all workloads is 1%.



**Figure 7: Performance improvements.**

### 5.3 Energy Savings

As expected, the EWT design achieved significant energy savings in writes. Figure 8 shows the write energy in each workload, normalized to the baseline. With EWT, up to 80% of write energy reduction is observed. Among all 17 workloads, 14 of them get more than 60% of write energy reduction. Even for workloads with lower bit write redundancy such as *mgrid*, *sphinx3* and *swim*, EWT still achieves 40%-60% savings. The average saving on write energy is 70%.



**Figure 8: Write energy savings.**

**Figure 9: Dynamic energy savings.**

We combined write energy and read energy together to evaluate our savings in total dynamic energy. Figure 9 shows the measured results. As write energy contributes to more than 70% of total dynamic energy in baseline, applying EWT leads to significant reduction in total dynamic energy (52% on average). We then compared the total energy (dynamic plus leakage) between EWT and baseline in Figure 10. With EWT, total energy can be reduced by up to 53%, and the average reduction is 33%.

### 5.4 Energy-Delay Product

Last, combining our results in both energy and performance, we present results for  $ED^2$  in Figure 11. Due to the significant savings in energy and slight improvement in execution time, we obtain up to 59% of  $ED^2$  reduction, with an average reduction of 34%. These results show that EWT can effectively improve energy efficiency using a STT-RAM cache.

## 6. PRIOR ART

There have been active efforts recently on STT-RAM designs. Most device-level studies focus on improving the MTJ properties, cell array structures, and prototyping. To name a few, Hosomi *et al.* fabricated a 4K bit STT-RAM using their tailored MTJ design in 0.18 $\mu$ m technology [3]. The goal was to demonstrate that STT-RAM is a prominent candidate for next generation memory due its high speed, low power and high scalability. Kawahara *et al.*

