Assuring Security and Reliability of Emerging Non-Volatile Memories

Mohammad Nasim Imtiaz Khan School of EECS Pennsylvania State University University Park, USA muk392@psu.edu

Swaroop Ghosh School of EECS Pennsylvania State University University Park USA szg212@psu.edu

Abstract— At the end of Silicon roadmap, keeping the leakage power in tolerable limit has become one of the biggest challenges. Several promising Non-Volatile Memories (NVMs) offering highdensity, high speed, and competitive reliability/endurance while eliminating leakage issues are being investigated. On one hand, the above-desired properties make emerging NVM suitable candidates to assist or replace conventional memories in memory hierarchy as well as to infuse compute capability to eliminate Von-Neumann bottleneck. On the other hand, their unique features such as high and asymmetric read/write current and persistence bring new threats to data security while compute-capability imposes new fundamentally different security challenges. Some of these memories are already deployed in full systems and as discrete chips. Therefore, it is utmost important to investigate the security issues of NVMs spanning the application space. This work makes pioneering contributions to this challenge through a holistic approach- from devices to circuits and systems using a combination of design and test methodologies to develop secure and resilient NVMs. The proposed attacks and countermeasures are validated on test boards using commercial NVM chips. Finally, this research has been tied to education by converting the test boards to design a modular and reproducible self-learning cybersecurity kit which has been piloted to train graduate and undergraduate students and K-12 teachers.

Keywords— NVM, security, reliability, test, countermeasures

I. INTRODUCTION

Conventional volatile memories such as, Static RAM (SRAM) and Dynamic RAM (DRAM) suffer from significant leakage power whereas conventional storage class Non-Volatile Memories (NVMs) (e.g. Flash) suffer from higher write energy, poor performance, and endurance. However, emerging NVMs (note, henceforth the term 'NVM' is used to denote 'emerging NVM') are beneficial due to their desirable properties - zero leakage, high-density, scalability, high endurance and CMOS compatibility [1]. Some examples of NVMs are Spin-Transfer Torque RAM (STTRAM) [2], Magnetic RAM (MRAM) [3], Phase Change Memory (PCM) [4], Resistive RAM (RRAM) [5] and Ferroelectric RAM (FeRAM) [6]. Due to promising aspects, NVMs are already being commercialized by industries e.g., Everspin [7], Adesto [8], Intel/Micron (PCM) [9] and Cypress [10]. Intel's 3D Xpoint memory [9] is a recent example of NVM's adoption as a cache for Solid State Drives (SSDs).

Unique characteristic of NVMs brings fundamentally unique

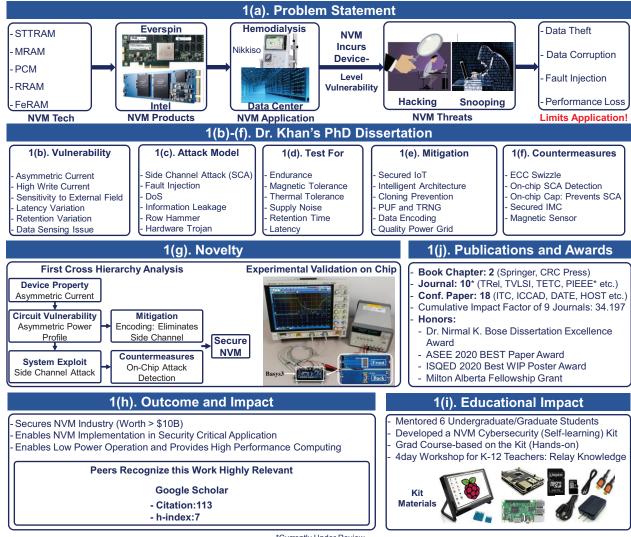
security issues that were not perceived before. Asymmetric read/write current (reveals data signature) and high write current (creates high supply noise) are examples to name a few. These vulnerabilities can be leveraged to launch Side Channel Attack (SCA) or fault injection attack (details in Section II). Interestingly, NVM can be used for storage such as cache, main memory as well as for In-Memory Computing (IMC), each bringing their own suite of security issues. For example, cache may employ 1T-1S (S: storage), main memory/storage may use 1D-1S (D: diode), and, IMC may use either of them but with conceptually different read/write/computing modes. Therefore, the same attack model may not be applicable to all application modes even for the same vulnerability. Similarly, a common mitigation technique might not be sufficient across all application modes. For example, advanced wear-leveling techniques (with high computational overhead) might be suitable for the main memory but not for cache due to tighter performance requirements. Therefore, NVM vulnerabilities, data security, and reliability issues should be investigated across application modes for a deeper understanding of the challenges and for the development of strong countermeasures (Fig. 1(a)).

Novelty of this study (Fig. 1(g)): This research is one of the firsts to perform a comprehensive study on NVM vulnerabilities (Fig. 1(b)) by investigating them from device, circuit and system point of view (for cache/storage/IMC application), and, proposes various attack models (Fig. 1(c)) to exploit these vulnerabilities and mitigations/countermeasures to secure them (Fig. 1(e)-(f)). This work also performs experimental validation of the attack models and countermeasures on commercial chips. Note that reliability and security issues are intertwined with each other since reliability issues can be exploited to launch attacks. This work also proposes testing techniques to detect NVM reliability issues that can become security vulnerabilities (Fig. 1(d)) in short test time/energy.

Relationship with prior works: Various NVM security issues e.g., cold-boot attack on PCM-based main memory [11-12] and SCA on MRAM [13], are investigated before. PCM-based main memory is also investigated for its susceptibility to birthday paradox attack [14] and to reduced-lifetime due to various workloads that can target write operations to specific cells [15-19]. Therefore, countermeasures such as Online Attack Detector (OAD) [16] and Dual Counter Encryption technique, namely DEUCE [19] are proposed. However, significant vulnerabilities and attack vectors remain unexplored.

1

TTTC-PhD INTERNATIONAL TEST CONFERENCE 978-1-7281-9113-3/20/ \$31.00 ©2020 IEEE



*Currently Under Review

Fig. 1 (a) Pictoral representation of the problem statement; (b)-(f) Dr. Khan's Ph.D. dissertation [26] showing, (b) NVM vulnerabilities, (c) proposed attack models, (d) proposed tests to identify vulnerability/reliability issues, (e) mitigation and (f) countermeasures proposed to secure the NVM; (g) the major novelty, (h) the high level outcome/impact; (i) the educational impact; (j)-(h) quality of Dr. Khan's thesis is measured through the number of publications, citations, honors and awards. The google scholar citation record in Fig. 1(h) shows that peers recognize this work highly relevant.

Conventional tests techniques either fail to capture NVM-specific issues or incur high test time [20]. Moreover, limited studies are performed on NVM tests. For example, [21-22] proposes tests to identify retention time while [23-24] proposes algorithms to capture various defects/coupling faults. However, several reliability issues (e.g. endurance, supply noise, etc.) need to be studied and a full test flow needs to be developed for NVM.

Impacts: The impact of this study are described below:

Technical impacts (Fig. 1(h)): This research will, i) secure NVM industry which is worth more than \$10B, ii) enable safe usage of discreet/on-chip NVMs in a broad range of devices including security-critical sectors e.g., healthcare and banking; and, iii) enable low-power operation for Internet-of-Things (IoTs) and High-Performance Computing (HPC) by removing

Von-Neuman bottleneck.

Educational impacts (Fig. 1(i)): During the course of the investigation, Dr. Khan has mentored and trained 6 graduate/ undergraduate students on designing test boards and characterizing security issues of commercial MRAM, RRAM, and FeRAM chips. Dr. Khan has repurposed these test boards to develop a self-learning kit, and design a hands-on curriculum on NVM cybersecurity to train graduate/undergraduate students and K-12 teachers. The learnings from the training kit has been selected as the Best Paper in the American Society of Engineering Education (ASEE) conference [25] which is the flagship venue for engineering education (Fig. 1(h)).

Rest of the paper is organized as follows: Section II introduces NVM and describes their vulnerabilities identified in

this work; **Section III** summarizes the proposed attack models; **Section IV** describes the test techniques proposed to identify vulnerabilities; **Section V** summarizes the mitigations/ countermeasures proposed in this work; **Section VI** describes the developed cybersecurity kit; **Section VII** draws conclusion.

II. NVMs, Their Applications, and Vulnerabilities

This section briefly introduces the NVMs and, covers their applications and vulnerabilities identified in this work. Although, we have considered five flavors of NVMs, for the sake of brevity, we summarize this paper taking STTRAM, MRAM, and RRAM as examples.

A. STTRAM and MRAM and Their Vulnerabilities

STTRAM cell (Fig. 2(a)) contains one Magnetic Tunnel Junction (MTJ) as the storage element which contains a free (FL) and a pinned (PL) magnetic layer. The resistance of the MTJ stack is high (low) if FL magnetic orientation is antiparallel (parallel) compared to the PL. MTJ can be toggled from parallel (P) (data '0') to anti-parallel (AP) (data '1') (or vice versa) using current induced Spin-Transfer Torque by passing the appropriate write current (> critical current) from source-line to bit-line (or vice versa). MRAM is like STTRAM except its write operation is electric field-driven. A current is passed through the MTJ top/bottom metal plates with appropriate polarity that creates a magnetic field to switch FL polarization.

Vulnerabilities: STTRAM/MRAM suffers from,

- Asymmetric write and read current (Fig. 2(c)-(d)): The asymmetric current drawn by the memory can be leveraged as data signature, and SCA can be launched [27];
- High write current: The high NVM write current (Fig. 2(c)) can lead to high supply noise such as supply voltage droop and ground bounce. This is true since the power grid is implemented at a higher metal layer (say M₈), and the bitcells are implemented at a lower metal layer (say M₁) (Fig. 2(e)). Therefore, there exists a significant parasitic capacitance and resistance between them, and the bitcells incur a supply voltage droop or ground bounce. This can be leveraged to launch fault injection [28], DoS [28], information leakage [29], and Row Hammer (RH) [30] attacks on a shared memory space;
- Susceptibility to external fields: External magnetic field can flip the magnetic orientation of MTJ FL, and corrupt the data [31]. This can be leveraged to launch DoS attack;
- Susceptibility to temperature: High temperature can lead to reduced data retention, and adversary can launch DoS;

• Sensitivity to Process Variation (PV): High PV results in weaker cells (low retention) which are vulnerable to attacks.

B. RRAM and its Vulnerabilities

RRAM contains an oxide material between its Top/Bottom Electrode (TE/BE) (Fig. 2(b)). RRAM's resistive switching is due to oxide breakdown and re-oxidation which modifies a Conduction Filament (CF). Conduction through the CF is primarily due to transportation of electrons in the oxygen vacancies. These vacancies are created under the influence of an electric field due to the applied voltage. The two states of the RRAM are termed as Low Resistance State (LRS) and High Resistance State (HRS). The process of switching the state to LRS (HRS) is known as SET (RESET).

Vulnerabilities: RRAM also suffers from, (i) high write current; (ii) asymmetric write and read current; (iii) susceptibility to temperature; and, (iv) low endurance. Malicious workload can hammer a RRAM cell to exhaust its life.

C. Vulnerabilities of PCM and FeRAM

PCM suffers from, (i) high write current; (ii) asymmetric write/read current; and, (iii) low endurance. Malicious workload can hammer a cell to exhaust its life. FeRAM suffers from, (i) high write current, (ii) asymmetric read current; and, (iii) susceptibility to external electric/thermal field which can flip the polarization of FeRAM ferroelectric material. Therefore, DoS attack can be launched by applying external fields.

D. Application-specific Vulnerabilities

• Vulnerabilities for memory application: NVMs can be integrated in various levels of memory hierarchy [2-6], and its vulnerabilities can vary based on the hierarchy since their specification is different. L1 cache has strict performance requirements (<0.5ns), uses physical addressing (to avoid address translation for high performance) and resides inside the core. Furthermore, the patterns in address and data signals are rare which presents unique scope to the adversary in inserting novel NVM hardware Trojan triggers and payloads (details in Section III.F). Such attacks are not possible for Last Level Cache (LLC)or main memory since logical to physical address translation naturally obfuscates the physical address to prevent adversary from hammering a target address to sensitize the trigger. Furthermore, LLC is slower (<1ns) and on-chip. The performance for main memory (~200ns) and SSD (10-100µs) are relaxed than on-chip cache. Therefore, data encryption is a feasible countermeasure (especially for SSD) but not suitable

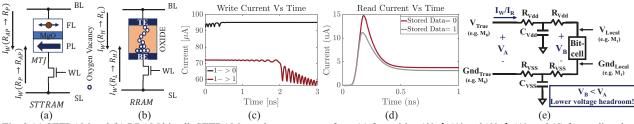


Fig. 2 (a) STTRAM and (b) RRAM bitcell; STTRAM supply current waveform (c) for writing '0' \rightarrow '1' and '0' \rightarrow '1', and (d) for reading data '0' and '1'. A significant gap is present between the write currents of '0' \rightarrow '1' and '0' \rightarrow '1' as well as read current for data '0' and '1' which can be leveraged as data signature; (e) cartoon showing parasitics between upper metal layer (e.g. M_8) and lower metal layer (e.g. M_1).

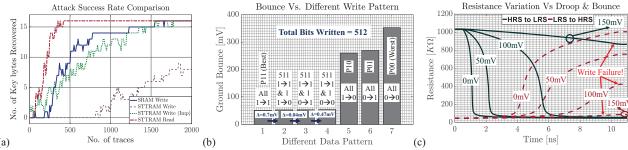


Fig. 3 (a) Comparison of the number of key bytes retrieved by SCA with respect to the number of current traces for SRAM and STTRAM; (b) supply noise (e.g., ground bounce) generated during write is dependent on write data; (c) RRAM write latency increases as supply noise increases.

for cache. Another example is SCA on L1 cache which might not be possible since it is deeply embedded in the chip, and requires very sophisticated measurement equipment. SCA on LLC and main memory are much less complex since the time window is large and signal processing can easily isolate other high frequency power components.

• Vulnerabilities for IMC application: NVM can bridge the widening performance gap between processor and memory. In Von-Neumann architecture, data must be fetched from memory to processor in order to compute. This adds a bottleneck to performance and energy-efficiency. NVMs have unique computation capability with faster operation and lower power compared to their volatile counterpart. Note that IMC using NVM requires read/write but with very different settings compared to normal memory operations. Therefore, the attack models that work on memory may not be applicable to IMC. For example, NVM-based IMC leaks computational signature and thereby, SCA on IMC architecture can lead to reverse engineering and reveal the underlying intellectual property.

III. ATTACK MODELS TO EXPLOIT NVM VULNERABILITES

This section summarizes the proposed data security and privacy attack models to exploit NVM vulnerabilities and steal/corrupt sensitive data.

A. SCA Leveraging Asymmetric Current [27]

SCA is a powerful threat, which targets weak implementation of cryptographic algorithms rather than the algorithm itself. The implementation weakness is related to the device physics of the underlying computing element, which makes it hard to fix the vulnerability. We have considered STTRAM as LLC and the Advanced Encryption Standard (AES)-128 with 10 rounds. The round outputs are stored and read from STTRAM LLC. An adversary has physical access to the system to monitor the power line, and can extract the read/write current. We have leveraged HD leakage model with Pearson correlation [32] and attacked the last round of AES to investigate SCA vulnerability.

Fig. 3(a) summarizes the result with the number of key bytes retrieved with respect to the required current traces for STTRAM and SRAM. The first byte of the key can be retrieved from STTRAM write current in around 600 traces which is suboptimal. The work further improved the attack model with some basic pre-processing (subtracting the average initial write

current from the final write current). Similar analysis is also performed on STTRAM read and SRAM write operation to compare the results. Before the pre-processing, STTRAM write was revealing 8 bytes of the key in around 2000 traces whereas after pre-processing, STTRAM write reveals all 16 bytes in around 1600 traces. Note that SRAM write also reveals 16 bytes in similar number of traces. However, STTRAM read operation is more vulnerable since it leaks all 16 bytes in just 400 traces.

We have verified this analysis on the read operation of a commercial MRAM chip which has a similar read operation to STTRAM. We have identified that a window of 15ns i.e., Window of Interest (WOI) performs sensing of actual data from the memory cells. Results show that the average read current depends on the number of ones in 8-bit read data. Next, we have implemented the SCA on MRAM read operation and successfully extracted the full key.

B. Fault Injection/DoS Attack Leveraging Supply Noise [28]

As mentioned in **Section II.A**, NVM write current is high and bitcells can incur high supply noise due to the parasitics between true V_{dd} /gnd and local V_{dd} /gnd (**Fig. 2(e)**). The magnitude of the noise depends write data pattern since I_{write} for $0 \rightarrow 0$, $0 \rightarrow 1$, $1 \rightarrow 0$ and $1 \rightarrow 1$ are different. **Fig. 3(b)** shows the noise generated by a full cache line write for various write data patterns.

RRAM write operation is simulated with additional supply noise generated by a parallel operation in another independent memory bank. As the bitcell being written incur more noise from the parallel operation, write latencies increases. Fig. 3(c) shows the RRAM resistance switching during write with respect to supply noise. Supply noise > 50mV and < 120mV can cause LRS to HRS write failure but still can write HRS to LRS. If the adversary can generate supply noise in a way that the victim incurs noise in this range, it can inject a 0→1 polarity fault. However, if the victim incurs supply noise > 120mV, it can cause complete write failure i.e. DoS attack.

C. Row Hammer Attack Leveraging Supply Noise [30]

If an adversary can keep writing to a particular address, the generated high ground bounce can propagate to the word-line/source-line/bit-line drivers of the neighboring bits. If the bounce propagates to word-lines drivers, the access transistors of the unselected bits sharing the same bit-line/source-line drivers with the selected cells, will partially turn ON and incur a disturb current. These bitcells may experience retention

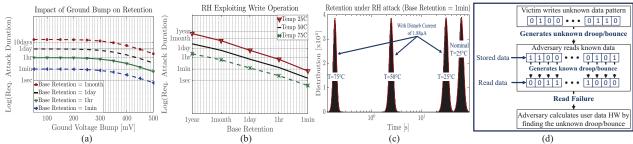


Fig. 4 (a) Impact of ground bounce on retention time of unselected bits (base retention = 1month); (b) impact of RH attack on STTRAM write for various base retentions; and, (c) retention distribution under RH attack (base retention = 1min). 1-million-point Monte-Carlo analysis is conducted with 3σ of 2% of MTJ thermal stability factor, Δ_0 (mean = 24.85, corresponding retention ~1min); (d) overview of information leakage attack.

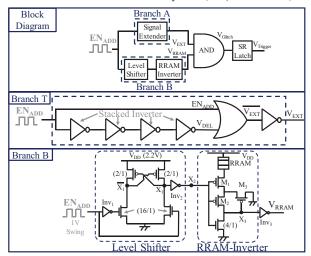


Fig. 5 One flavor of the dealy-based NVM Trojan trigger [34].

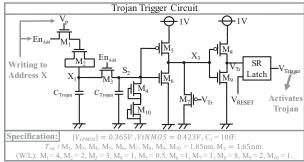


Fig. 6 Capacitor-based Trojan trigger circuit [35].

failure and read disturb. Furthermore, if the bounce propagates to source-line/bit-line drivers, the bitcells will experience lower voltage headroom and read/write operations to them may fail.

Fig. 4(a) shows that as the ground bounce seen by the bitcell increases, the retention time of the cell reduces. A higher temperature can reduce the retention time further (**Fig. 4(b)**). The RH attack can flip the bits in \sim 30secs at T = 25°C if the base retention is 1 min which can be reduced to 2.5secs and 0.2secs at T = 50°C and T = 75°C, respectively. **Fig. 4(c)** shows the weaker bits under PV are more vulnerable to the attack since their retention reduces to \sim 19secs, \sim 1.7secs, and \sim 0.1secs at T

= 25°C, T = 50°C and T = 75°C, respectively.

D. Information Leakage Attack Leveraging Supply Noise [29]

Fig. 4(d) shows the concept of an information leakage attack by leveraging supply noise. Victim writes sensitive data patterns which creates data-dependent supply noise and propagates to the adversary memory space. Adversary reads a known data (i.e. known supply noise) which adds up to the propagated noise and creates a read failure. From the read failure characteristics, the adversary can detect the amount of sensitive noise from the victim and back-calculate the Hamming Weight of the victim's sensitive data. Further details of the attack can be found in [29].

E. DoS Attack by External Field [30]

STTRAM is susceptible to contactless tampering efforts, i.e. by subjecting it to a strong external magnetic field and/or thermal field, an adversary can corrupt the stored contents [31]. The PL of MTJ in STTRAM is robust. However, the FL of MTJ could be toggled through both spin-polarized current as well as a magnetic field. The FL is susceptible to both the magnitude and polarity of the external magnetic field it is subjected to. Note that all spintronic memories are expected to experience a similar issue, and hence vulnerable to tampering. In [31], we have proposed sensors to detect the attack, stall operation, and recover the data through Error Correcting Code (ECC).

F. NVM and Capacitor-based Hardware Trojans [34-35]

In [34], a delay (Fig. 5) and voltage based NVM Trojan trigger is proposed by exploiting the RRAM resistance drift under pulsing current. Simulation results indicate that these triggers can be activated by accessing a pre-selected L1 cache address 2500-3000 times (varies with trigger designs). The proposed trigger evades the testing phase since it requires a high number of hammerings. It can also evade system-level detection techniques that can classify hammering as a potential security threat since hammering needs not be consecutive due to non-volatility. Once triggered, this Trojan can launch several attacks [34]. The area/static/dynamic power overheads of the trigger are $6.68 \mu m^2/104.24 \mu W/0.426 \mu W$, respectively.

In [35], a capacitor-based Trojan trigger (Fig. 6) for NVM is proposed which is small, sneaky, and stealthy. A pre-defined L-1 cache address can be hammered with a predefined data pattern. Every hammering increases the charge stored in a capacitor, and when it is charged more than a threshold, it

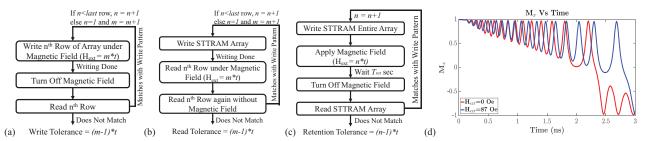


Fig. 7 Algorithm for finding, (b) write, (c) read, and (c) retention tolerances; (d) free layer magnetic orientation (M_x) vs time for write tolerance testing. Noted that 86Oe magnetic field can be tolerated during write (@87Oe, M_x is \sim -0.75 i.e. the bitcell is not written successfully) [36].

generates a signal (considered as Trojan trigger). It can evade detection during testing since it requires a large number of hammering. Optical inspection may also not work since the circuit is small. The circuit also consumes low static power which makes it difficult to detect via power spectrum comparison with a golden chip. However, the limitation of such a capacitor-based Trojan is the hammering requires to be fairly continuous. If hammering is stopped for a sufficiently long period of time, the capacitor may get completely discharged.

IV. PROPOSED TESTS TO IDENTIFY NVM VULNERABILITY

This section describes the test techniques proposed to identify NVM vulnerability with short test time and cost.

A. Magnetic Tolerance Test [36]

Spintronic memories should be tested and certified during the write, read and retention mode separately since their tolerance for different modes could be different. Furthermore, chip to chip tolerance can identify the weakest chip due to PV. Therefore, a tolerance test could discard chip which has a lower tolerance than the rated one. If a chip still incurs an attack more than the threshold value, a sensor can detect it, and take necessary measures as proposed in [31].

Write Tolerance: Write tolerance is the maximum magnetic field under which it can be written successfully at a specified write current with a specified write latency.

Read Tolerance: Read tolerance is the maximum magnetic field under which it can be read successfully without causing any disturbance to the bits at a specific read current with a specified read latency.

Retention Tolerance: Retention tolerance is the maximum external magnetic field under which it does not incur any data-corruption for a specified time period during retention mode.

The algorithms to find the write tolerance, read tolerance and retention tolerance are shown in Fig. 7(a), 7(b) and 7(c), respectively [36]. Fig. 7(d) shows an example of write tolerance identification. During write operation with an external magnetic field of 87Oe, free layer magnetic orientation (M_x) is \sim -0.75 after 3ns. MTJ may or may not go back to its previous state $(M_x = 1)$ which indicates that the bit is not written successfully. Therefore, it is considered as a write error (for magnetic field > 86Oe). Write tolerance depends on the data that is being written whereas, read and retention tolerance depends on the stored data. The worst-case write tolerance occurs when writing $0\rightarrow 1$ since it requires higher write time.

The worst-case read and retention tolerance of a bit occur when the bit stores data '1' (since data '0' (P state) is the preferred state for STTRAM i.e. writing $1 \rightarrow 0$ requires less energy).

B. Thermal Tolerance Test [37]

Memory chips should be certified to operate successfully within a target temperature range (typically: -10°C to 90°C):

- (a) At high temperature, the energy barrier between two states of the memory reduces. Therefore, the data retention time reduces and can lead to retention failure. Furthermore, read failure can occur since the reduction of resistance difference between two states leads to a reduction of sense margin and read disturb can occur since slight disturbance can flip the data at a lower energy barrier. Therefore, the manufacturer needs to test retention and read failure/disturb at high temperature.
- (b) At low temperature, the energy barrier between the two states increases. Therefore, the read/write latencies increase and can lead to read/write failures.

Thermal tolerance test can be done using the algorithms proposed for the magnetic tolerance test (Section IV.A) by applying an external thermal field instead of the magnetic field. This test can be combined with standard hot-cold test. The highest temperature at which the memory read failure/disturb does not occur and retention time meets the target specification is the upper limit of thermal tolerance. The lowest temperature at which read/write operation does not fail is the lower limit.

C. Supply Noise Testing [20, 38]

An adversary can leverage supply noise to launch various attacks as discussed in **Section III.B-D**. This is especially true if parallel operations are done to near by independent memory banks and they have a strong coupling (lower resistance and capacitance between two addresses of two independent banks). Mostly, the victim bitcells are the weaker cells due to PV.

In [20], a supply noise test technique is proposed to capture the impact of parallel write operations on the weaker bitcells. The method is further improved in [38] which implements test time compression technique leveraging unique data patterns. The work also divides the testing scenarios into 3 cases as shown in **Fig. 8**: i) accesses in adjacent banks (e.g. Bank₀-Bank₁); (ii) accesses in physically confronting banks (e.g. Bank₀-Bank₂); and (iii) accesses in diagonal banks (e.g. Bank₀-Bank₃). These cases successfully captures a weaker bitcells which gets affected by near-by parallel operations. Techniques such as Wordline Overdrive (WLOV) and Early Write are

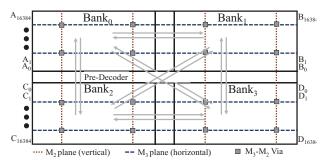


Fig. 8 4MB LLC diagram showing addresses of $Bank_0/Bank_1/Bank_2/Bank_3$ as A/B/C/D from 0 to 16K. The grey arrows show different test cases for testing the impact of supply noise. [38]

proposed to reduce write latency during testing which in turn reduces test time/energy. Once identified, either parallel accesses can be restricted to such weak cells or the chips with weaker bitcells can be discarded.

D. Endurance Testing [37]

NVM performance can degrade over time due to physical breakdown (STTRAM/RRAM/PCM) or resistance drift (RRAM/PCM). STTRAM/MRAM have an oxide layer in their storage element, MTJ, and RRAM has oxide layer between two electrodes in its bitcell. Oxide might breakdown due to high $I_{\rm write}$ leading to function failure. Note that LRS changes by 2X-10X and HRS changes by 5X-100X in TaO2 based RRAM due to variation. In PCM, time-dependent resistance drift in amorphous chalcogenide material is one of the major reliability concerns. Therefore, RH attack on NVM to target cells and exhaust its lifetime can be a big security concern.

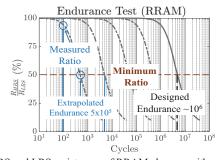


Fig. 9 HRS and LRS resistances of RRAM changes with respect to its usage cycle. This can be leveraged to design endurance test [37].

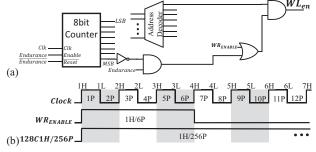


Fig. 10 (a) Proposed DFT circuit for endurance test; (b) input waveforms of the proposed DFT circuit (Fig. 10(a)) [37].

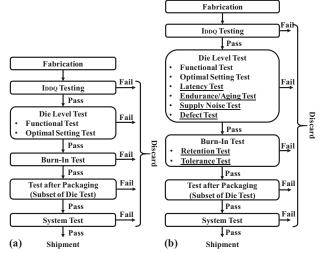


Fig. 11 (a) Conventional memory (e.g. SRAM/DRAM) test flow; (b) repurposed test flow for NVM [37].

We have proposed a novel test technique which can measure the endurance of the memory cell in a very short test time. The basic idea is to create the model of the physical parameters that changes as the cell is written multiple times. For example, Fig. 9 shows the change of ratio of resistance in HRS and LRS of RRAM bitcell with respect to the number of times they are written. Once such modeling is done, the bitcells of a chip can be hammered using a DFT circuit (Fig. 10). The bitcell resistance can be measured, and the effective endurance can be calculated by leveraging the relation shown in Fig. 9 and applying extrapolation. If the endurance is lower than the target threshold, they can be considered as vulnerable chips.

E. Retention and Latency Tests [37, 39]

In [39], we have proposed Magnetic Burn-In (MBI) for retention testing for spintronic memory. The idea is to apply external magnetic field antiparallel to the magnetic orientation of MTJ FL. This will reduce retention time of the cell and in turn, the overall test time. Combining MBI with conventional thermal Burn-In (BI) (applying both external magnetic and high thermal field), namely MBI+BI further reduces the test time. In [37], we have also proposed latency test which captures the long latency tail accurately. Additionally, we proposed WLOV and Early Write techniques to reduce the test time by 2.15X.

F. Summary [37]

Fig. 11(a) shows the test flow for conventional memories and Fig. 11(b) shows the repurposed test flow proposed in this work for NVM. Table I summarizes all the NVM test methods proposed in this work along with prior works stating the required time for these tests and the test time/energy compression achieved by this work. The total test time found for a 4MB NVM is around 5.05 secs considering all the tests with available test times. Note that all tests are not required for each memory type. We approximate that on an average 3/4 secs are required to test each chip which is close to typical test time (2/3sec) considered [40] for each chip. Therefore, the proposed

Test Name	Application	Proposed Methods (Our work is shown in bold reference)	Test Time (Compression)	Test Energy Compression
Retention Test	All NVMs	Weak write-based [21]	13.3s (N/A)	-
		EMACS [37]	15.6s (N/A)	-
		External Temperature-based [32]	-	-
		MBI [39]	$3.42s (1.68 \times 10^6 X)$	-
		MBI+BI [39]	$0.25s (4.12 \times 10^7 X)$	
Latency Test	All NVMs	WL OV + Early Write [37]	87ms (2.15X)	4.97µJ
Supply Noise Test	All NVMs	[20]	646.23s (N/A)	•
		WL OV + Early Write [38]	1.57s (410.82X)	79.88J
Magnetic Tolerance Test	Spintronic Memories	[36]	262.14ms (write) (N/A)	
			327.68ms (read) (N/A)	-
			458.72ms (retention)(N/A)	
Thermal Tolerance Test	All NVMs	[37]	524.28ms (write) (5X)	
			655.36ms (read) (5X)	-
			917.44ms (retention) (5X)	
Endurance Test	All NVMs	[37]	$50.02 \mu s (9.99 \times 10^3 X)$	$9.99 \times 10^{3} X$

TABLE I SUMMARY OF VARIOUS PROPOSED NVM TESTS [37]

tests are effective in keeping the test time close to the target.

V. MITIGATIONS AND COUNTERMEASURES TO SECURE NVM

This section summarizes the mitigation techniques and the countermeasures proposed in this work to secure NVM.

A. Mitigations

- 1) Secured IoT with STTRAM Replacing eFlash [31]: We have investigated the challenges to replace eFlash in IoTs with STTRAM. We have considered a non-invasive magnetic field attack and assumed that only one IoT in a homogenous network is under attack at a time. This situation is likely when multiple IoTs are distributed in a building or critical infrastructure such as a bridge, to collect the required information. During the normal operation, if an adversary tries to attack the STTRAM to corrupt the Firmware (FW), the attack sensors and the Integrity Checker is able to sense the attack ahead of time. The Integrity Checker sends the HALT interrupt to microcontroller after attack detection. If the attack is launched when the IoT is powered off, the passive sensors can still detect the attack due to the failure of sensor bits. When the IoT is powered up after the attack, the boot ROM triggers the STTRAM Integrity Checker and the STTRAM will fail the integrity check due to the modified sensor arrays from the previous attack. The boot ROM then starts executing the recovery code from the EPROM.
- 2) Prevent Accessing Same Address Repeatedly [30]: Write operation can be stalled to facilitate recovery of lost retention to mitigate the susceptibility of STTRAM to RH attack. The average disturb current reduces by 80% by stalling write by one cycle after every four consecutive writes which in turn increases the attack duration by $1.30 \, \mathrm{X}$ and $1.57 \, \mathrm{X}$ for 1min of base retention at $T = 50 \, \mathrm{^{\circ}C}$ and $T = 75 \, \mathrm{^{\circ}C}$, respectively.
- 3) Duplicate Write Driver [30]: Duplicate column write drivers can be incorporated to cut down write current. This in turn reduces supply noise by half. If noise reduces from 300mV to 150mV, the RH attack duration will increase from 37.69secs to 1min (retention unaffected) for 1min of base retention.

- 4) Good Quality Power Grid [30]: A 25% reduction of paracitic resistance in the power grid reduces noise by 25%.
- 5) Data Encoding [27]: Data encoding can obfuscate data signature. We have investigated several types of data encoding, and found that they increase the attack complexity.
- 6) Intelligent Architecture: We have shown that if the system is intelligent enough to detect two parallel operations going to memory arrays that are physically separated by a minimum resistance (23 Ω for our design), and allow only one operation at a time, supply noise induced-attacks such as fault injection, information leakage and RH attacks can be prevented. However, this introduces a performance overhead. We have also proposed memory controller that can implement various advanced techniques such as 'Read before write', 'Changing logical to physical mapping after a threshold' etc. to improve wear-leveling and get more life out of NVM.
- 7) Extensive March Test at Scaled Supply Voltage [34]: This can accelerate Trojan triggering and can detect the payloads being executed by identifying power profile anomaly.

B. Countermeasures

- 1) On-Chip SCA Detection Sensor [41]: We have leveraged the dependency of Ring Oscillator frequency on the supply voltage and proposed an on-chip sensor that can successfully detect SCA when adversary inserts a resistance in the power rail. This technique incurs minimal area overhead and is resilient to process and temperature variation.
- 2) ECC Swizzle: To prevent an attacker from adding malicious Trojans to tamper with the ECC bits, we propose scrambling of the ECC columns between different address rows using hardware bit swizzling. Adversary will fail to inject faults using hardware Trojan since ECC bits cannot be tampered with. If a chip incurs frequent ECC failures in particular addresses, that can be considered as an attack, and the system can change logical to physical mapping to avoid accessing them.

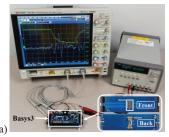






Fig. 12 (a) CSTM01 and (b) CSTM02 module; (c) K-12 teacher workshop to train them on using the cybersecurity kit [25].

- *3) On-Chip Capacitor-based SCA Prevention:* We have proposed on-chip capacitor that isolates NVM from the supply during the read/write operations. This prevents adversary from detecting the power drawn by NVM during read/write operations.
- *4) Secure IMC:* We have proposed redundant inputs and expansion of literals to secure NVM-based IMC from SCA with minimal overheads. These proposed techniques obfuscate the computation signature of the underlying IP of IMC preventing reverse engineering.

C. Morphable Security Primitive [42]

We have proposed morphable security primitive that can be used both as Physically Unclonable Function (PUF) and True Random Number Generator (TRNG) [42] by manipulating the write time and the number of write pulses to a commercial toggle MRAM chip. PUF is used for authentication and can prevent counterfeiting and reverse engineering of NVM products. Furthermore, TRNG is used in cryptographic applications and can secure data stored in NVM. The proposed security primitives are based on the observation that intrinsic and extrinsic variations in the MRAM MTJ changes its write latency. Therefore, if the write latency is chosen carefully, it will statistically toggle the same bit in two different chips (useful for PUF). Furthermore, the same bit will be stochastically flipped if written multiple times with different write time (useful for TRNG).

VI. CYBERSECURITY KIT ON NVM SECURITY [25]

This work develops a hands-on and modular self-learning Cybersecurity Training (CST) kit to advance cybersecurity education. Students can promptly apply newly acquired knowledge on the kit. This kit accompanies Do-It-Yourself training modules that is used to model and investigate cybersecurity issues and their prevention to all levels of the cybersecurity workforce, including undergraduate and graduate students and K-12 science teachers. In this section, two modules (out of four) of the kit are described which are on NVM security.

A. CSTM01: FPGA Board for Side Channel Attacks on NVM

Modern cryptographic algorithms deployed in current systems are mathematically sound and secure. However, due to the specific nature of implementation in hardware, they may leave some side-channel signatures that can be leveraged to attack the cryptosystem. We have prepared a module with a FPGA (Basys 3), MRAM, breadboard and SCA resistance as

shown in Fig. 12(a), that can be used to exploit a hardware implementation of a cryptosystem. Software along with the hardware are provided to investigate the security issues of several cryptosystems such as AES, RSA etc. The code is flexible enough to be modified by the students to experiment.

B. CSTM02: FPGA Board for Magnetic Attacks on NVMs

STTRAM and MRAM are susceptible to external magnetic field. An adversary can corrupt the stored data in NVM by applying a higher external magnetic field compared to its tolerance. We have prepared a module (Fig. 12(b)) using a Xilinx Virtex-5 FPGA which is programmed to run specific March tests on the Device Under Test (DUT). The FPGA is interfaced with a PC to acquire test results for flexible test automation and data analysis. This approach allows for easy replacement of the DUT and avoids any need for manual entry of results. We have used commercially available toggle MRAM. Furthermore, this framework is modular which is amenable to new test routines and other flavors of NVM.

A graduate-level hands-on course has been introduced in the School of EECS, Penn State on hardware security. A 4-day workshop (Fig. 12(c)) has also been organized to train K-12 science teachers on cybersecurity aspects of computing who in turn plan to pilot the training kit in their respective classes.

VII. CONCLUSION

We have identified fundamental NVM vulnerabilities, proposed new attack models, and finally, proposed mitigations and countermeasures to provide data security and privacy guarantees. Several studies are also validated experimentally on commercial chips. We have proposed test techniques to capture vulnerabilities in short test time/energy. Finally, a self-learning modular kit is designed to disseminate the knowledge to graduate and undergraduate students. This work can enable NVM implementation in security critical applications and secure NVM industry which is worth more than \$10B.

ACKNOWLEDGMENT

This work is supported by NSF (CNS - 1722557, CCF - 1718474, DGE - 1723687 and DGE - 1821766), DARPA Young Faculty Award (D15AP00089) and SRC task 2847.001. We also thank our collaborators A. Chattopadhyay, S. Bhasin, S. Gupta, J. Park, R. Jha, S. Thirumala, A. Jones, A. Yuan, A. Iyengar, H. Motaman, A. De, A. Saki, K. Nagarajan, S. Sayyah and N. Gattu.

REFERENCES

- Mark H. Kryder and Chang Soo Kim. "After hard drives—What comes next?." Magnetics, IEEE Transactions on 45, no. 10 (2009): 3406-3413.
- [2] Xiuyuan Bi, Zhenyu Sun, Hai Li, and Wenqing Wu. "Probabilistic design methodology to improve run-time stability and performance of stt-ram caches." In Proceedings of the International Conference on Computer-Aided Design, pp. 88-94. ACM, 2012.
- [3] Jing Li, Patrick Ndai, Ashish Goel, Sayeef Salahuddin, and Kaushik Roy. "Design paradigm for robust spin-torque transfer magnetic RAM (STT MRAM) from circuit/architecture perspective." Very Large Scale Integration (VLSI) Systems, IEEE Transactions on 18, 2010.
- [4] Zhou, Ping, Bo Zhao, Jun Yang, and Youtao Zhang. "A durable and energy efficient main memory using phase change memory technology." In ACM SIGARCH Computer Architecture News, ACM, 2009.
- [5] Stuart Schechter, Gabriel H. Loh, Karin Straus, and Doug Burger. "Use ECP, not ECC, for hard failures in resistive memories." In ACM SIGARCH Computer Architecture News, ACM, 2010.
- [6] Y. M. Kang and S. Y. Lee, "The challenges and directions for the massproduction of highly-reliable, high-density 1T1C FRAM,"17th IEEE International Symposium on the Applications of Ferroelectrics, 2008.
- [7] MR4A08BUYS45 Data Sheet, [Online], Available www.everspin.com/file/882/download, [Accessed: Jul 16, 2020].
- [8] CBRAM Technology, [Online], Available: www.dialogsemiconductor.com/products/memory/cbram-technology, [Accessed: Jul 16, 2020].
- [9] Intel Optane, [Online], Available: ark.intel.com/content/www/us/en/ark/products/series/190349/inteloptane-persistent-memory.html, [Accessed: Jul 16, 2020].
- [10] FM28V102A, [Online], Available: www.cypress.com/file/140901/download, [Accessed: Jul 1, 2020].
- [11] S. Chhabra and Y. Solihin "i-NVMM: a secure non-volatile main memory system with incremental encryption", 2011 38th Annual International Symposium on Computer Architecture (ISCA), San Jose, CA, 2011.
- [12] P. Zhou, B. Zhao, J. Yang and Y. Zhang⁺, "A durable and energy efficient main memory using phase change memory technology", SIGARCH Comput. Archit. News, 2009.
- [13] A. Chakraborty, A. Mondal and A. Srivastava, "Correlation power analysis attack against STT-MRAM based cyptosystems," 2017 IEEE International Symposium on Hardware Oriented Security and Trust.
- [14] André Seznec, "Towards Phase Change Memory as a Secure Main Memory." in IEEE Computer Architecture Letters, Jan. 2010.
- [15] Jingfei Kong and Huiyang Zhou, "Improving privacy and lifetime of PCM-based main memory," 2010 IEEE/IFIP International Conference on Dependable Systems & Networks (DSN), Chicago, IL, 2010.
- [16] Vinson Young, Prashant J. Nair, and Moinuddin K. Qureshi, "DEUCE: Write-Efficient Encryption for Non-Volatile Memories", In Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '15).
- [17] Elnawawy, Hussein, Mohammad Alshboul, James Tuck and Yan Solihin. "Efficient Checkpointing of Loop-Based Codes for Non-volatile Main Memory." 2017 26th International Conference on Parallel Architectures and Compilation Techniques (PACT) (2017): 318-329.
- [18] Moinuddin Qureshi, John Karidis, Michele Franceschini, et. al, "Enhancing lifetime and security of PCM-based main memory with start-gap wear leveling". 2009 42nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO), New York, NY, 2009.
- [19] Moinuddin K. Qureshi, Andre Seznec, Luis A. Lastras, and Michele M. Franceschini, "Practical and secure PCM systems by online detection of malicious write streams", 2011 IEEE 17th International Symposium on High Performance Computer Architecture (HPCA '11).
- [20] Mohammad Nasim Imtiaz Khan and Swaroop Ghosh, "Test challenges and solutions for emerging non-volatile memories," 2018 IEEE 36th VLSI Test Symposium (VTS), San Francisco, CA, 2018, pp. 1-6.
- [21] A. Iyengar, S. Ghosh and S. Srinivasan, "Retention Testing Methodology for STTRAM," in IEEE Design & Test, Oct. 2016.
- [22] I. Yoon, A. Chintaluri and A. Raychowdhury, "EMACS: Efficient MBIST architecture for test and characterization of STT-MRAM

- arrays," 2016 IEEE International Test Conference (ITC), 2016.
- [23] A. Chintaluri, A. Parihar, S. Natarajan, H. Naeimi and A. Raychowdhury, "A Model Study of Defects and Faults in Embedded Spin Transfer Torque (STT) MRAM Arrays," 2015 IEEE 24th Asian Test Symposium (ATS).
- [24] A. Chintaluri, A. Parihar, S. Natarajan, H. Naeimi and A. Raychowdhury, "A Model Study of Defects and Faults in Embedded Spin Transfer Torque (STT) MRAM Arrays," IEEE 24th Asian Test Symposium (ATS), 2015.
- [25] Asmit De, Mohammad Nasim Imtiaz Khan, Karthikeyan Nagarajan, Abdullah Ash Saki, et al. "Hands-on Cybersecurity Curriculum Using a Modular Training Kit". 2020 ASEE Virtual Annual Conference.
- [26] Mohammad Nasim Imtiaz Khan, "Assuring Security and Privacy of Emerging Non-Volatile Memories", The Pennsylvania State University, Ann Arbor, ProQuest Dissertations Publishing, 2019. 28097059.
- [27] Mohammad Nasim Imtiaz Khan, Shivam Bhasin, Alex Yuan, Anupam Chattopadhyay and Swaroop Ghosh, "Side-Channel Attack on STTRAM Based Cache for Cryptographic Application," 2017 IEEE International Conference on Computer Design (ICCD), Boston, MA, 2017, pp. 33-40.
- [28] Mohammad Nasim Imtiaz Khan and Swaroop Ghosh, "Fault injection attacks on emerging non-volatile memory and countermeasures", In Proceedings of the 7th International Workshop on Hardware and Architectural Support for Security and Privacy (HASP '18), 2018.
- [29] Mohammad Nasim Imtiaz Khan and Swaroop Ghosh, "Information Leakage Attacks on Emerging Non-Volatile Memory and Countermeasures", In Proceedings of the International Symposium on Low Power Electronics and Design (ISLPED '18), 2018.
- [30] Mohammad Nasim Imtiaz Khan and Swaroop Ghosh, "Analysis of Row Hammer Attack on STTRAM," 2018 IEEE 36th International Conference on Computer Design (ICCD), Orlando, FL, USA, 2018, pp. 75-82.
- [31] Asmit De, Mohammad Nasim Imtiaz Khan, Jongsun Park and Swaroop Ghosh, "Replacing e-Flash with STTRAM In IoTs: Security Challenges And Solutions", Journal of Hardware and Systems Security, Dec. 2017.
- [32] Paul Kocher, Joshua Jaffe, and Benjamin Jun, "Differential Power Analysis", In AICC, Springer Berlin Heidelberg, 1999.
- [33] Swaroop Ghosh, Mohammad Nasim Imtiaz Khan, Asmit De and Jae-Won Jang, "Security and privacy threats to on-chip Non-Volatile Memories and countermeasures," 2016 IEEE/ACM International Conference on Computer-Aided Design (ICCAD), Austin, TX, 2016, pp. 1-6.
- [34] Karthik Nagarajan, Mohammad Nasim Imtiaz Khan, and Swaroop Ghosh, "ENTT: A Family of Emerging NVM-based Trojan Triggers", IEEE International Symposium on Hardware Oriented Security and Trust (HOST), 2019.
- [35] Mohammad Nasim Imtiaz Khan, Karthik Nagarajan, and Swaroop Ghosh, "Hardware Trojans in Emerging Non-Volatile Memories", Design, Automation & Test in Europe Conference & Exhibition (DATE), 2019.
- [36] Mohammad Nasim Imtiaz Khan, Anirudh S Iyenga and Swaroop Ghosh, "Novel Magnetic Burn-In for Retention and Magnetic Tolerance Testing of STTRAM," in IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 26, no. 8, pp. 1508-1517, Aug. 2018.
- [37] Mohammad Nasim Imtiaz Khan and Swaroop Ghosh, "Test Methodologies, and, Test Time Analysis and Compression for Emerging Non-Volatile Memory," In IEEE International on Reliability, 2019.
- [38] Mohammad Nasim Imtiaz Khan and Swaroop Ghosh, "Test of Supply Noise for Emerging Non-Volatile Memory," 2018 IEEE International Test Conference (ITC), Phoenix, AZ, USA, 2018, pp. 1-10.
- [39] Mohammad Nasim Imtiaz Khan, Anirudh S Iyenga and Swaroop Ghosh, "Novel magnetic burn-in for retention testing of STTRAM," Design, Automation & Test in Europe Conference & Exhibition (DATE), 2017.
- [40] Addressing Test Time Challenges, [Online], Available: semiengineering.com/addressing-test-time-challenges/, [Accessed: Jul 16, 2020].
- [41] Navyata Gattu, Mohammad Nasim Imtiaz Khan, Asmit De and Swaroop Ghosh, "Power side channel attack analysis and detection," 2020 International Conference on Control, Automation and Diagnosis, 2020.
- [42] Mohammad Nasim Imtiaz Khan, Chak Yuen Cheng, Sung Hao Lin, et. al., "A Morphable Physically Unclonable Function and True Random Number Generator using a Commercial Magnetic Memory", 2020 21st International Symposium on Quality Electronic Design (ISQED), 2020.