# Reliable Low-Power High Performance Spintronic Memories

For obtaining the academic degree of

## Doctor of Engineering

Department of Informatics
Karlsruhe Institute of Technology (KIT)
Karlsruhe, Germany

**Approved**

## Dissertation

by

## Rajendra Kumar Bishnoi

From India

| | |
|---|---|
| Date of Oral Examination: | 17.05.2017 |
| Adviser: | Prof. Dr. Mehdi Baradaran Tahoori, KIT, Germany |
| Co-adviser: | Dr. Guillaume Prenat, Spintec, France |

Rajendra Kumar Bishnoi
Berliner-ring 2,
76676 Graben-neudorf

Hiermit erkläre ich an Eides statt, dass ich die von mir vorgelegte Arbeit selbstständig verfasst habe, dass ich die verwendeten Quellen, Internet-Quellen und Hilfsmittel vollständig angegeben haben und dass ich die Stellen der Arbeit - einschließlich Tabellen, Karten und Abbildungen - die anderen Werken oder dem Internet im Wortlaut oder dem Sinn nach entnommen sind, auf jeden Fall unter Angabe der Quelle als Entlehnung kenntlich gemacht habe.

Karlsruhe, Mai 2017

_____

Rajendra Kumar Bishnoi

# ABSTRACT

On the pathway of Moore's law, chip manufacturing industries have shown explosive growth for the last five decades. With that, the demand of memory components have also increased exponentially, leading to a memory dominant chip in today's computing system. However, the traditional on-chip memory technologies such as *Static Random Access Memories* (SRAMs), *Dynamic Random Access Memories* (DRAMs) and flip-flops are facing severe challenges in terms of scaling, leakage power and reliability. These challenges along with the overwhelming demand of increase in performance and density for an on-chip storage, make researchers seek for a new non-volatile storage technology. Emerging spintronic based storing technologies such as *Spin Orbit Torque* (SOT) and *Spin Transfer Torque* (STT) have captured a lot of attention in recent years because of their various beneficial features such as non-volatility, scalability, high endurance, CMOS compatibility and soft-error immunity. In the spintronic technology, spin of electrons represent the information, in which the data is stored as resistance states, that can be altered by passing a polarized current through the storing device. The zero leakage capabilities of its storing devices as well as attributes of normally-off/instant-on computing, make it very effective to deal with static power challenge. However, this technology faces some challenges, such as high access latency and energy consumption, before its widespread utilization. Addressing these challenges requires new computing paradigms, architectures and design philosophy.

In the spintronic technology, the issue of high access latency is because the storing cell takes relatively high switching durations, more than conventional SRAM cell. Moreover, due to stochastic switching nature of the storing cell and the influence of process variation, a significant timing margin is required. For that period of time, a constant write current has to flow through the bit-cell to ensure its switching, resulting in a huge energy consumption. Equivalently, the read operation also requires a considerable delay margin because of the impact of process variation. On the other hand, due to extra timing margins, read and write currents flow for a longer duration, leading to several other reliability problems such as read disturb and other degradation issues such as *Time Dependent Dielectric Breakdown* (TDDB). Hence, it is necessary to reduce energy and latency as well as to improve reliability, in order to make it a potential candidate for an on-chip storage system.

In this thesis work, we provide design solutions to address the challenges associated with caches, register and flip-flop designs using a cross-layer approach. For caches, we designed several circuit-level techniques, whose efficiencies are evaluated at both memory architecture-level and system-level in terms of energy consumption, performance improvement and reliability enhancements. We design a *self-timed technique* for both read and write operations for caches which can dynamically detect the completion of their respective operations. The read and write operations are terminated immediately after this detection in order to save energy. Moreover, the self-timed technique shortens the current flow durations from the storage devices, that improve TDDB and read disturb effects during the write and read operations, respectively. To improve write latency, we boost the write current through the bit-cell to accelerate the magnetic switching process. On the top of that, to address register-file related challenges, we design a multi-port memory architecture in which we have exploited a unique feature of the SOT cell that it can perform simultaneous read and write operations. In this way, read-write contention can be resolved at the bit-cell level, resulting in a much simplified multi-port register file architecture.

In addition to memory solutions, we have also proposed two flip-flop architectures. The first one is a *Non-Volatile Non-Shadow* flip-flop architecture, in which storing devices are employed

as active components. This allows immediate turn-off/turn-on the supply voltage, hence it is beneficial for aggressive power gating. Overall, the proposed flip-flop design has similar timing characteristics as conventional CMOS flip-flops for normal operations, and at the same time it allows to reduce the static power significantly compared to shadow non-volatile flip-flops. The second one is a *Fault Tolerate* flip-flop architecture which is resilient to various defects and faults. The effectiveness of all proposed techniques is illustrated with extensive circuit-level simulations, which are further supported by detailed system-level evaluations. Overall, using our proposed techniques, we have achieved significant improvements for performance, energy and reliability for spintronic on-chip storages such as caches, registers, flip-flops and latches.

# ZUSAMMENFASSUNG

Moores Gesetz folgend, ist es der Chipindustrie in den letzten fünf Jahrzehnten gelungen, ein explosionsartiges Wachstum zu erreichen. Dies hatte ebenso einen exponentiellen Anstieg der Nachfrage von Speicherkomponenten zur Folge, was wiederum zu speicherlastigen Chips in den heutigen Computersystemen führt. Allerdings stellen traditionelle on-Chip Speichertechnologien wie *Static Random Access Memories* (SRAMs), *Dynamic Random Access Memories* (DRAMs) und Flip-Flops eine Herausforderung in Bezug auf Skalierbarkeit, Verlustleistung und Zuverlässigkeit dar. Eben jene Herausforderungen und die überwältigende Nachfrage nach höherer Performanz und Integrationsdichte des on-Chip Speichers motivieren Forscher, nach neuen nichtflüchtigen Speichertechnologien zu suchen. Aufkommende spintronische Speichertechnologien wie *Spin Orbit Torque* (SOT) und *Spin Transfer Torque* (STT) erhielten in den letzten Jahren eine hohe Aufmerksamkeit, da sie eine Reihe an Vorteilen bieten. Dazu gehören Nichtflüchtigkeit, Skalierbarkeit, hohe Beständigkeit, CMOS Kompatibilität und Unanfälligkeit gegenüber Soft-Errors. In der Spintronik repräsentiert der Spin eines Elektrons dessen Information. Das Datum wird durch die Höhe des Widerstandes gespeichert, welche sich durch das Anlegen eines polarisierten Stroms an das Speichermedium verändern lässt. Das Problem der statischen Leistung gehen die Speichergeräte sowohl durch deren verlustleistungsfreie Eigenschaft, als auch durch ihr Standard- Aus/Sofort-Ein Verhalten an. Nichtsdestotrotz sind noch andere Probleme, wie die hohe Zugriffslatenz und die Energieaufnahme zu lösen, bevor sie eine verbreitete Anwendung finden können. Um diesen Problemen gerecht zu werden, sind neue Computerparadigmen, -architekturen und -entwurfsphilosophien notwendig.

Die hohe Zugriffslatenz der Spintroniktechnologie ist auf eine vergleichsweise lange Schaltdauer zurückzuführen, welche die von konventionellem SRAM übersteigt. Des Weiteren ist auf Grund des stochastischen Schaltvorgangs der Speicherzelle und des Einflusses der Prozessvariation ein nicht zu vernachlässigender Zeitraum dafür erforderlich. In diesem Zeitraum wird ein konstanter Schreibstrom durch die Bitzelle geleitet, um den Schaltvorgang zu gewährleisten. Dieser Vorgang verursacht eine hohe Energieaufnahme. Für die Leseoperation wird gleichermaßen ein beachtliches Zeitfenster benötigt, ebenfalls bedingt durch den Einfluss der Prozessvariation. Dem gegenüber stehen diverse Zuverlässigkeitsprobleme. Dazu gehören unter Anderem die Leseintereferenz und andere Degenerationspobleme, wie das des *Time Dependent Dielectric Breakdowns* (TDDB). Diese Zuverlässigkeitsprobleme sind wiederum auf die benötigten längeren Schaltzeiten zurückzuführen, welche in der Folge auch einen über längere Zeit anliegenden Lese- bzw. Schreibstrom implizieren. Es ist daher notwendig, sowohl die Energie, als auch die Latenz zur Steigerung der Zuverlässigkeit zu reduzieren, um daraus einen potenziellen Kandidaten für ein on-Chip Speichersystem zu machen.

In dieser Dissertation werden wir Entwurfsstrategien vorstellen, welche das Ziel verfolgen, die Herausforderungen des Cache-, Register- und Flip-Flop-Entwurfs anzugehen. Dies erreichen wir unter Zuhilfenahme eines Cross-Layer Ansatzes. Für Caches entwickelten wir verschiedene Ansätze auf Schaltkreisebene, welche sowohl auf der Speicherarchitekturebene, als auch auf der Systemebene in Bezug auf Energieaufnahme, Performanzsteigerung und Zuverlässigkeitverbesserung evaluiert werden. Wir entwickeln eine Selbstabschalttechnik, sowohl für die Lese-, als auch die Schreiboperation von Caches. Diese ist in der Lage, den Abschluss der entsprechenden Operation dynamisch zu ermitteln. Nachdem der Abschluss erkannt wurde, wird die Lese- bzw. Schreiboperation sofort gestoppt, um Energie zu sparen. Zusätzlich limitiert die Selbstabschalttechnik die Dauer des Stromflusses durch die Speicherzelle, was wiederum das Auftreten von TDDB und Leseinterferenz bei Schreib- bzw. Leseoperationen reduziert. Zur Verbesserung der Schreiblatenz heben wir den Schreibstrom an der Bitzelle an, um

den magnetischen Schaltprozess zu beschleunigen. Um registerbankspezifische Anforderungen zu berücksichtigen, haben wir zusätzlich eine Multiport-Speicherarchitektur entworfen, welche eine einzigartige Eigenschaft der SOT-Zelle ausnutzt, um simultan Lese- und Schreiboperationen auszuführen. Es ist daher möglich Lese/Schreib- Konfilkte auf Bitzellen-Ebene zu lösen, was sich wiederum in einer sehr viel einfacheren Multiport- Registerbankarchitektur niederschlägt.

Zusätzlich zu den Speicheransätzen haben wir ebenfalls zwei Flip-Flop-Architekturen vorgestellt. Die erste ist eine nichtflüchtige non-Shadow Flip-Flop-Architektur, welche die Speicherzelle als aktive Komponente nutzt. Dies ermöglicht das sofortige An- und Ausschalten der Versorgungsspannung und ist daher besonders gut für aggressives Powergating geeignet. Alles in Allem zeigt der vorgestellte Flip-Flop-Entwurf eine ähnliche Timing-Charakteristik wie die konventioneller CMOS Flip-Flops auf. Jedoch erlaubt er zur selben Zeit eine signifikante Reduktion der statischen Leistungsaufnahme im Vergleich zu nichtflüchtigen Shadow- Flip-Flops. Die zweite ist eine fehlertolerante Flip-Flop-Architektur, welche sich unanfällig gegenüber diversen Defekten und Fehlern verhält. Die Leistungsfähigkeit aller vorgestellten Techniken wird durch ausführliche Simulationen auf Schaltkreisebene verdeutlicht, welche weiter durch detaillierte Evaluationen auf Systemebene untermauert werden. Im Allgemeinen konnten wir verschiedene Techniken entwickeln, die erhebliche Verbesserungen in Bezug auf Performanz, Energie und Zuverlässigkeit von spintronischen on-Chip Speichern, wie Caches, Register und Flip-Flops erreichen.

# Contents

# Glossary

A | C | D | E | F | G | I | L | M | N | O | P | R | S | T | V | W

## A

**AP-Magnetization**

   Anti-Parallel Magnetization, when the magnetic orientation of two ferromagnetic layers of the MTJ cell are in the opposite direction.

## C

**CACTI**

   An integrated cache and memory access time, cycle time, area, leakage, and dynamic power model.

**CLK**

   Clock signal of the design.

**CMOS**

   Complementary Metal Oxide Semiconductor.

## D

**DRAM**

   Dynamic Random Access Memory.

## E

**ECC**

   Error Correction Code.

## F

**FabScalar**

   A heterogeneous multi-core processor with configurable out-of-order instruction.

**FF**

   Flip-Flop.

**FRAM**

   Ferroelectric Random Access Memory.

**FTNV-L**

   Fault Tolerant Non-Volatile Latch.

## G

*Glossary*

**Gem5**

Cycle acutate performance simulator.

**GND**

Ground Supply Voltage.

**I**

**ISCAS**

Sequential Benchmark Circuits.

**ITRS**

International Technology Roadmap for Semiconductors.

**L**

**Leon3**

CPU microprocessor core which has an in-order seven stage pipeline unit.

**M**

**MiBench**

Commercially representative embedded benchmark suite for various applications.

**MRAM**

Magnetic Random Access Memories.

**MTJ**

Magnetic Tunnel Junction.

**N**

**NMOS**

N-Channel Metal Oxide Semiconductor Field Effect Transistor.

**NVM**

Non-Volatile Memory.

**NVNS-FF**

Non-Volatile Non-Shadow flip-flop.

**NVSim**

A circuit-level module for non-volatile memory performance, energy and area estimation.

**O**

**OpenSPARC**

Processor with an in-order, six stage pipeline that executes four threads concurrently.

**P**

**P-Magnetization**

Parallel Magnetization, when the magnetic orientation of two ferromagnetic layers of the MTJ cell are in the same direction.

**PCRAM**

Phase-Change Random Access Memory.

**PG**

Power Gating.

**PMOS**

P-Channel Metal Oxide Semiconductor Field Effect Transistor.

**R**

**RA**

Product of Resistance and Area of the MTJ cell.

**RAM**

Random Access Memory.

**RRAM**

Resistive Random Access Memory.

**S**

**SoC**

System-on-Chip.

**SOT**

Spin Orbit Torque.

**SPEC**

Standard Performance Evaluation Corporation CPU benchmark.

**Spectre**

Cadence tool for circuit simulations.

**SPICE**

Simulation Program with Integrated Circuit Emphasis.

**SRAM**

Static Random Access Memory.

**SRPG**

State Retention Power Gating.

**STT**

Spin Transfer Torque.

**T**

**TDDB**

Time Dependent Dielectric Breakdown.

**TMR**

Tunnel Magneto-Resistance.

**V**

**VDD**
    Power Supply Voltage.

**W**

**WE**
    Write Enable.

**WER**
    Write Error Rate.

# List of Figures

# List of Tables

# 1 Introduction

In this era of the information age, electronic equipments heavily flood the market in almost every application fields such as communication, household appliances, automotive, emerging internet of things, education, research, medical, security, etc. The computing systems play a significant role in the development of such equipments [26–29]. In other words, the development of new products or upgrading the existing ones, can only become possible because of the continuous progress of computing systems. For that, the main credit goes to the semiconductor industries for their explosive growth, which they have shown in the last five decades on the pathway of Moore's law. Moore's law is a projection made by Gordon E. Moore in 1965 [30], which states that *the number of transistors on integrated circuits doubles approximately every two years.* As per Moore's law, not only a high device density on a chip is achieved, but also the devices become faster, consume less energy and cheaper every year [31–33]. For instance, the *Skylake* microprocessor, one of the recent Intel's processor, has around 1.75 billion transistors and CPU clock frequency of 4 GHz, which is around five orders of magnitude denser and approximately three orders of magnitude faster compared to the early days microprocessors [4]. With advancements in the computing system, the demand of memory components have also grown exponentially [29].

In general, in a microprocessor, memories are organised hierarchically in order to balance the memory response time and their capacity requirements. A typical memory hierarchy consists of memory mass storages such as disk and flash, main memory, various level of caches, registers and flip-flops. The memory mass storages are the class of *secondary memory*, that are well known for their high capacity, nevertheless these memories are extremely slow. On the other hand, the main memory, which is typically built using *Dynamic Random Access Memories* (DRAM), is relatively faster than the mass storage system, however, much slower compared to the speed of the processor. Additionally, the rate of improvements in DRAM is slower over the years compared to that of processor, which results in processor speed improvements to be



Figure 1.1: Conventional memory trends (a) Trends of DRAM capacity, data obtained from [1–3], (b) Trends of total SRAM-caches capacity, data obtained from [4].

Figure 1.2: Dynamic power verses static power for various technology nodes (data obtained from [5]). Total energy consumption is increased with scaling in which static power is dominating.

masked by that of DRAM speed, well-known with the term *memory-wall* in a microprocessor [34, 35]. Therefore, in order to meet performance requirements of the processor, multiple layers of buffers in the form of caches are introduced between the main memory and processor unit. Cache memory, which is developed using *Static Random Access Memories* (SRAM) technology, is faster, but less dense and more expensive than DRAM. Moreover, the other on-chip storage systems such as registers and flip-flops, which are also implemented using SRAMs, are much smaller and faster than caches. The DRAM as a primary memory and SRAM as an on-chip memory, demonstrated tremendous growth in the past to fulfill the requirement of the computing system. Figure 1.1(a) and Figure 1.1(b) show the rate of increase of memory capacity per year for DRAM and SRAM, respectively. With these, memories nowadays have become a dominating component of the modern computing system. For instance, in certain microprocessors, the embedded memories are already occupying more than 70 % of the total die area [36], a trend that is increasing as technology advances.

With this technology advancements, the energy consumption at the chip-level has also increased exponentially, already reaching at the alarming-level, that is well-known as the *power-wall* in a microprocessor, and memory components contribute mostly to it. The reason is that, first, memories are dominating the total chip components and second, memories are not only consuming power when they are actually operating but also when they are not in operational mode. This is due to the fact that both SRAM and DRAM are volatile memories, which means their storage cells always require a constant power supply to retain their value. Therefore, SRAM based memories significantly contribute to static (leakage) power consumption, and DRAM cells require a periodical refresh. As a consequence, nowadays, the leakage power has even exceeded the operational mode of power, and the trend is that this ratio is increasing further with the technology down-scaling.

With the shrinking technology dimensions, SRAM and DRAM technologies are also facing several other issues due to *scalability* and reliability [37–39]. A direct impact of the scaling is that the storing cell itself becomes unstable, meaning the storing capability of the cell becomes weak. As per *International Technology Roadmap for Semiconductors* (ITRS) guidelines, device dimension reduction would be practically impossible after 2025 [40]. This is because, the *metal-oxide-semiconductor* transistors have already reached a countable dopant atoms. Further scaling results in a significant variations in the device parameters, which in turn widely varies the device strength of the small bit-cells. As a consequence, the design constraints such as read,

Figure 1.3: Spintronic based storing devices.

write and hold margins can be severally impacted. Furthermore, the smaller storing devices become more vulnerable to soft-errors due to radiations and external noises.

In these memories, the fundamental storage principles are based on the charge of the electrons. However, the other characteristic of electrons, that is, the spin of electrons, is completely ignored. Researchers have explored this quantum property of the electron as well as its associated magnetism to come up with an emerging technology named as *spintronic*. Spintronics means *spin transfer electrons* (or *spin based electrons*) is the new science of computer and memory chips that are based on electron spin [41, 42]. For memory solutions, this technology is known as *Magnetic Random Access Memories* (MRAM) such as *Spin Transfer Torque* (STT) and *Spin Orbit Torque* (SOT), in which the storing device is *Magnetic tunnelling Junction* (MTJ) cell which consists of two ferromagnetic layers, separated by a thin oxide layer. These MTJ cells are employed to store the digital information, and the magnetic orientation of these devices are manipulated to store 1's or 0's in the form of the resistance states. These cells are not only non-volatile in nature, i.e. not losing their content after being turned-off (means zero leakage), but also proven to be highly scalable. In addition, these devices have several other beneficial features such as high density, the device can sustain virtually infinite number of write cycles (i.e., high endurance), compatible to CMOS, immune to soft-error due to radiations, etc [43]. Overall, this technology has all the beneficial features as the existing memory technologies like speed of SRAM, density of DRAM and non-volatility of Flash. Moreover, this technology has an edge over the other emerging technologies such as *Phase-change*, *Resistive* and *Ferro-electric* memories, in terms of scalability, access latencies and endurance. Hence, spintronic-based memories have potential to become a universal memory technology, that can be fit in any memory hierarchy of the computer system.

Despite all these advantages, spintronic based memories have relatively high write access latency and energy consumption in comparison to SRAM memories [44]. These are due to the physical characteristic of the storing device as described next. Firstly, in this technology, the switching latency of the MTJ cell itself is high as it requires a high write current for a long duration (several nanoseconds) to switch their magnetization. Secondly, the write latency of the bit is asymmetric in nature which means different switching delay for write '1' and write '0' [45–48]. Thirdly, the switching of the MTJ cell is stochastic in nature, which means switching latency of the cell is random even for the same environmental and operating conditions [25, 49]. Hence, to attain a required *Write Error Rate* (WER), a significant timing margin is required to account the stochastic effects, which is in additional to the margin required for process variations. This timing margin is added with respect to the one which has the worst switching delay, i.e. slow-writes. Due to these timing margins, the total write period increases significantly. For that period of time, a constant write current has to flow through a MTJ cell

Figure 1.4: Illustration of the memory hierarchy, scope of the thesis work and contributions.

to ensure its switching completion, resulting in a huge overall energy consumption. On the top of that, several reliability challenges are also associated with spintronic based memories. For instance, due to the non-deterministic switching nature of the cell, there is a possibility that the required magnetic transition does not happen within the given write duration, leading to a write error. On the other hand, due to extra timing margins, read and write currents flow for a longer duration, generating other reliability problems. For instance, a larger read period increases the possibility of unintended switching in a MTJ, known as *read disturb*, whereas a high write current flows for a longer duration can accelerate the degradation of MTJ parameters as well as lifetime, termed as *Time Dependent Dielectric Breakdown* (TDDB). *The goal of this thesis is to improve overall performance, energy and reliability in spintronic based memories using cost effective solutions.*

## 1.1 Thesis Scope and Contributions

In this thesis work, we propose design techniques to address the challenges associated with spintronic based caches, registers and flip-flops, using a cross-layer approach as shown in Figure 1.4. In that, we exploited the cell-level properties of this technology using circuit design and evaluated their performance and energy consumption using a system-level platform. The details of our contributions are described as:

**Contributions for on-chip Caches**

- **Self-timed memory write operations**: Spintronic based memories have a long write period due to the requirement of a significant timing margin to account stochastic switching, however, the write operations for most of bits will finish much before the actual write termination time. Such cells consume unnecessary energy even after the completion of their switching because their switching completion time is unknown. In order to generate a per-bit write completion signal, we have developed a flip detection circuit. The write operation for each bit-cell is terminated based on this signal. Using this technique, a significant amount of energy is reduced as here energy is consumed only for the durations when the memory write components are actually performed. Moreover, for bit-cells which are already in their required magnetization state, the write completion signal is generated at the beginning of the operation, hence energy savings are even more in such cases. Additionally, using this self-timed memory write technique, we can detect the write errors as well as the impact of TDDB and read disturb can be reduced significantly as the duration of the current flow is reduced.

4

- **Write speed acceleration technique**: Using the aforementioned self-timed write scheme, the write energy can be reduced significantly, however, the total write period still remains the same to target a certain WER value. To full-fill the performance requirement of the on-chip caches, it is necessary to reduce the total write period along with the energy reduction of the spintronic-based memories. Therefore, the latency of the slow-write operation and the associated write margins have to be reduced in order to shorten the overall write period. In our proposed technique, we boost the current passing through the MTJ cell to accelerate the switching process. The current boosting technique works in two phases: 1) the write current is increased only for the slow-writes to the extent that its latency matches with that of the fast-writes. 2) We furthermore increase the write current in a step-wise manner only for the unfinished bit-cells. The proposed technique can significantly improve the write latency for a given WER because of the shortening of the timing margin. Moreover, due to the reduction in the total write period, the overall energy efficiency is increased as well as the effect of TDDB is reduced.

- **Self-timed memory read operations**: Similar to write, the read operation in spintronic-based memory requires a significant timing margin due to process variation. Due to this, the read current flows for a long time period, resulting in, not only a high energy consumption but also a high rate of read disturb failures. In general, the read operation is performed using a sense amplifier which always generates two complementary outputs. These two output signals are exploited dynamically in this proposed technique, when the read is performed, to generate a read detection signal. Once read detection signals for all bit-cells are generated, a combined acknowledgement signal is generated to deactivate the read operation in order to save energy. The self-timed read technique can significantly reduce the energy consumption as well as the read disturb rate, because in this technique, the sense amplifier remains active only for the necessary duration.

### Contributions for Registers

- **Design of unique multi-ported memory architecture**: Multi-port memory is the key element to achieve massive parallel data processing demand in a microprocessor. With the continuous strive to increase the performance, the number of such memory elements (like registers) or/and their associated capacity grows significantly. As a consequence, the leakage is considerably increased because the state-of-the-art multi-ported memories use SRAM-based storing cells. Additionally, during read-write contention, i.e., when the same address is accessed simultaneously by a read and a write requests, a prioritizing or buffering approach is typically employed. This severely impacts the overall performance and also introduces a lot of complexity in the design. We design a unique multi-port memory architecture using a non-volatile SOT cell, which has no bit-cell leakage. Furthermore, we observe a fact about an SOT based storing cell that it can inherently perform simultaneous read and write operations. That means, both read and write operations can be performed by a single cell in the same cycle. In this way, read-write contention can be naturally resolved at the bit-cell level, resulting in a much simplified multi-processing design.

### Contributions for Flip-flop

- **Design of non-shadow flip-flop architecture**:

  A non-volatile spintronic-based flip-flop design is very beneficial to deal with static power challenge as unlike the conventional SRAM-based flip-flop designs, the power of the entire logic-core can be completely turned-off during standby durations. In these flip-flops, the

non-volatile components are employed as a shadow (or backup) latch, where the data backup storage can be done locally for each flip-flop. However, switching between backup and restore operations for such flip-flops requires a complex controlling mechanism which contributes to the delay, area and power overheads. Moreover, these flip-flops are only feasible for long sleep durations, making them ineffective for the short standby durations. We propose a novel non-volatile non-shadow flip-flop architecture using SOT as it offers shorter switching delays at lower current, which allows very aggressive power gating. Since MTJs are active components in this architecture, this architecture is suitable for short standby periods as well. For further energy efficiency, we utilize a redundant write avoidance scheme, i.e. if the value to be written is already stored in the MTJ cells, the write operation for those MTJ cells is not initiated. To study the impact of our proposed design for various applications, a system-level analysis is also performed. Overall, our proposed flip-flop architecture has similar timing characteristics as normal CMOS flip-flops in the active mode, and at the same time, it can address short standby periods unlike the shadow flip-flop architecture.

- **Design of fault-tolerant flip-flop architecture**: The spintronic-based storing device is more prone to manufacturing defects than the CMOS technology because these are fabricated based on a new process and new materials. In order to render manufacturing defects, memories are usually equipped with redundancies and error detection/correction mechanisms. However, these techniques are inapplicable to flip-flop design because flip-flops are scattered widely in the system-on-chip layout as individual cells. Therefore, a single fault in flip-flops can lead to the breakdown of the entire non-volatile latch scheme for a given design block. A traditional solution to resolve this issue is *Triple Modular Redundancy* for the non-volatile latch, nevertheless it incurs overall huge area, energy and performance costs. Therefore, we design a fault tolerant non-volatile flip-flop architecture, which is resilient to various MTJ related faults. It has almost the same performance and energy consumption as the standard non-volatile flip-flop design with a negligible area overhead. The proposed flip-flop architecture is applicable to both shadow as well as non-shadow flip-flop architectures.

## 1.2 Outline

The remainder of the thesis is primarily categorized into four parts. In the first part, Chapter-2 contains memory classification, necessary background and state-of-the-art about the spintronic technology. In the second part, Chapter-3 and Chapter-4 cover cache related proposed techniques to improve the overall energy consumption, performance and reliability. The third part, which is Chapter-5, focusses on the design of a novel multi-ported memory architecture for a register-file application. The fourth part is about the flip-flop architectures, in which Chapter-6 and Chapter-7 explain the design of a non-volatile non-shadow flip-flop and a fault tolerant flip-flop architecture, respectively. A brief introduction of each chapter is described as follows:

Chapter-2 describes first the technical background in general about semiconductor memories such as their classification, merits and demerits in a computing system, etc. Afterwards, the evaluation of spintronic technology in chronical order is described. Additionally, various switching mechanisms associated with the spintronic technology is mentioned in this chapter.Then, a cross layer design approach starting from bit-cell architecture all the way up to the entire memory architecture is described with their system-level evaluation in this chapter.

Chapter-3 proposes a self-timed read and write operation in MRAM technology. In this chapter, a circuit-level technique to generate a bit-wise completion signal for both write and read operation is described, followed by their system-level evaluations. This chapter also discuss

the reliability implications such as write errors, read disturb and cell degradation, after proposed self-timed techniques are employed.

In Chapter-4, the proposed performance improvement techniques are explained. First, a circuit-level static technique is explained, in which the latency of only slow-writes is reduced to nullify the effect of asymmetry in latency. Then, a circuit-level dynamic technique is described, in which the write current is increased in a step-wise manner only for those bits that have unfinished transitions. Finally, a system-level evaluation and cell degradation effects for our proposed techniques are explained in the end of this chapter.

Chapter-5 contains the description of a novel multi-port memory architecture using SOT-MRAM. In this chapter, first the implementation of our proposed design is explained, in which the bit-cell modification and all possible design scenarios with respect to our proposed architecture are discussed. Additionally, the two low power techniques, namely, *Unnecessary Write Termination* and *Unnecessary Write Avoidance*, that are introduced on the top of our proposed multi-port memory architecture, are described.

Chapter-6 describes a novel *Non-volatile Non-shadow* Flip-flop architecture in which storing cells are demonstrated as active components of the latching operation. Moreover, a write avoidance circuitry is designed for this flip-flop architecture to enhance its energy efficiency, that is explained latter in this chapter. Finally, a comprehensive comparison with CMOS-based flip-flop and their system-level evaluation are described at the end of the chapter.

In Chapter-7, a novel fault tolerant flip-flop architecture is explained. In the implementation part, a comprehensive analysis of its functionality in the presence of various defect is discussed. Later, an algorithm is mentioned, which is to determine the basic design parameters of the storing cell, so that the proposed design can generate the correct output in the presence of defects. In the end of the chapter, the functionality of the proposed architecture is analyzed in the presence of the temperature and process variations.

Finally, Chapter-8 concludes the thesis and discusses the potential future direction of the spintronic technology.

# 2 Background

Semiconductor memories have demonstrated a significant improvement and expansion in the past several decades, achieved higher speeds, lower energies, higher densities, more advanced features and lower costs. With that, the memory system has become one of the most critical components of any computing system. In this chapter, first a general description about various semiconductor memories along with their merits/demerits in a computing system is explained. Afterwards, the evaluation of spintronic technology is discussed. Then, characteristics of spintronic technology is explained. Afterwards, a bit-cell architecture and basic read/write circuits are explained, followed by a spintronic based flip-flop architecture using MTJ cells is mentioned. Finally, a memory architecture-level design parameter generation with their system-level evaluation is explained in the end of the chapter.

## 2.1 Classification of semiconductor memories

Memories are physical devices that are capable of storing information temporarily or permanently. They play a vital role in a computing system in terms of storing programs during their execution as well as data storage. In order to facilitate the requirement of application-oriented systems such as high-performance, energy-efficient, high-data intensive, low-cost, etc., different classes of memories are employed depending on their characteristics. Classification of memory array system is illustrated in Figure 2.1. Depending on applications, these memory arrays are broadly classified into three categories namely, *Random Access Memories* (RAM), *Serial Access Memory* (SAM) [50–52] and *Content Addressable Memory* (CAM) [53, 54]. The RAM is a class of memory devices that allows data item to be read or written in any order. This is in contrast

Figure 2.1: Memory classification (figure partially obtained from [6]). Here, * → indicates conventional memories.

Figure 2.2: Bit-cell architectures for conventional memories. Here, 6T→ six transistors; 1T1C→ one transistor and one capacitor.

to the SAM where data is accessed sequentially, that means no address is necessary for these memories. The shift registers and queues are the example of such type of memories as shown in the figure. On the other hand, CAM is a different type of memory that determines the address(es) for a given data using a search and match operation.

The RAM can be categorized into two part, namely, *volatile memories* and *non-volatile memories*. The volatile memories are the type of memories that can retain their content as long as a power supply is applied, whereas non-volatile memories can hold data even in the absence of the power supply. Volatile memories are further divided into static-RAM (SRAM) and dynamic-RAM (DRAM) depending on their cell characteristics [6, 55, 56]. Moreover, there are several type of memories that are non-volatile such as *Read Only Memory* (ROM) [6, 56, 57], *Erasable Programmable ROM* (EPROM) [6, 57, 58], *Electrically Erasable Programmable ROM* (EEPROM) [6, 57, 59], Flash [60–62] and emerging memories [63, 64]. The ROM memories are further divided into *Masked-ROM* and *Programmable-ROM*. In the masked-ROM, the stored values are hard-wired during fabrication, whereas Programmable-ROM can be programmed by blowing up fuses using a high programmable voltage after the fabrication [6]. On the other hand, both EPROM and EEPROM are re-programmable memories. The EPROMs can be erased by an exposure to ultraviolet light, whereas EEPROMs can be erased electronically. Flash memories [65] are the improved version of the EEPROM with the main advantage that the value can be erase or programmed at the block-level unlike the byte-level in EEPROM. There are two types of Flash memories namely, *NAND-* and *NOR-*flash. The NAND-flash is ideal for the high capacity storage because of its multi-level storage attributes, where as NOR is best suited for the storage as well as executions due to its fast read access capabilities.

In addition to the conventional memory technologies presented above, there are several emerging memory technologies such as *Magnetoresistive RAM* (MRAM) [66–68], *Phase-Change RAM* (PCRAM) [69–71], resistive RAM (RRAM) [72–74] and Ferroelectric RAM (FRAM) [64, 75]. All these emerging memories are the part of the non-volatile memory technology. MRAMs are the spintronic storing technology in which values are stored using spin property of the electron and the value is stored in terms of resistance state. PCRAMs and RRAMs are another resistance type memories which exploits the material property to change their resistance value. PCRAMs store data using sufficient heating effect, that can change amorphous to crystalline state or vice-versa, whereas in RRAM data is stored based on the motion of ions under the influence of an electric field.

The SRAM, DRAM and Flash memories are largely adopted by the semiconductor industries, those are termed as *conventional memories*. The bit-cell architectures for these conventional memories are demonstrated in Figure 2.2. A six transistors (6T) for SRAM, and a single transistor and a single capacitor (1T1C) for DRAM, are the most commonly used bit-cell architectures. Whereas, Flash-NAND uses a single floating gate transistor for the storage, and for *Multi-level cell* (MLC) design, the same bit-cell is capable to store multiple bits. In memories,

(a) General memory architecture  (b) Memory array organization

Figure 2.3: Illustration of a typical memory architecture and a hierarchical memory organization.

bit-cells are arranged as a regular structure known as *memory core* and additional circuitry that facilitate the required memory operations is known as *memory periphery*. A typical memory architecture is shown in Figure 2.3(a). In most of memories, the memory array is organised with one row per word[1] and one column per bit in each word. The address is first latched, and then corresponding *word-line* (WL) is activated using a decoder logic. Once the bit-cells are selected using WL activations, the required memory operation is performed using respective read/write circuits. For a high density memory, the memory array is organised in a hierarchical way as shown in Figure 2.3(b). Bank is the top level structure and each bank consists of multiple mat. Similarly, each mat consists of multiple subarrays and these subarrays are the basic building blocks of the array organization. Every subarray has its own set of periphery circuitry.

The basic storing cell for SAM and CAM memories are similar to that of the SRAM. Also, flip-flops and latches, that are known for a single bit storage, are prevalent design used in a system-on-chip, also employ SRAM-based storing mechanism. Overall, in a typical computing system, the SRAM-based storing systems are predominating in on-chip components, DRAMs are well known for the main memory and Flash memories are used in large as high storage system. These conventional memories have performed well in the past, however, as per ITRS report, now they have reached to the saturation level due to their scaling limit [76–79]. Additionally, the total energy consumption per chip, in which leakage is dominating, is a serious issue, especially for battery-operated portable devices. The introduction of new non-volatile spintronic based memory technologies, such as MRAMs, promise to bring a revolutionary advancement in the landscape of memory system [41]. The details about the MRAM technology are discussed in subsequent sections in this chapter.

## 2.2 Evolution of spintronic technology

Spintronics is a branch of science that uses the spin of electrons phenomenon for the information processing in a computing system. The spin of electrons can be either *spin-up* or *spin-down*, representing the two states of the binary data in spintronic in a similar way as the conventional electronic uses charge to represent the information as zeros and ones. Such spin transfer in solid-state devices is achieved by exploiting the magnetic properties of ferromagnetic materials. The physics behind today's spintronic technology has been known for a long time. For

---

[1]Word is a collection of bits that a particular processor can handle

Figure 2.4: Magnetoresistance observation for Fe and Cr structure for the applied magnetic field for the experiments performed by [7]. Here, $H_s$ is the field to overcome the antiferromagnetic coupling.

instance, the interdependency of the resistances due to the respective magnetic orientation and current was first observed in 1856 by William Thomson [80]. Furthermore, the spin dependent conduction was first observed by Mott in 1936, in which the resistive characteristics of ferromagnetic metals at the curie temperature was examined [81, 82]. Later, an experimental study about the spintronic based magnetic domain walls was performed by Berger in 1978 [83]. However, the widespread interest in magnetic nano-structures began with the discovery of *Giant Magnetoresistance* effect, which is also known as the foundation step for spintronic technology.

### 2.2.1 Giant Magnetoresistance

The *Giant Magnetoresistance* (or simply GMR) effect was first observed by Albert Fert and Peter Grünberg, independently in year 1988 [7, 84], based on that, later in 2007, they jointly awarded Nobel Prize in Physics. The GMR effect is a relative change in the resistance observed over the multi-layered structure of ferromagnetic and non-ferromagnetic materials. There are many experiments conducted on various set of materials to evaluate their GMR effects. A large magnetoresistance impact was first demonstrated on Fe/Cr metal layers [85, 86]. The Magnetoresistance curves for Fe/Cr multilayer structure with respect to magnetic field for the experiments conducted by [7] is shown in Figure 2.4. This experiment is conducted at 4.2 K, where Fe (Iron) layer thickness is fixed at 3 nm and Cr (Chromium) layer thickness is varying for 9 nm, 12 nm and 18 nm. Here, the magnetoresistance value was observed as high as 80 % for Fe-3 nm/Cr-0.9 nm multilayer structure. This experimental result drew a lot of attention in the field of nano-magnetism.

The GMR was also explored in other multilayer structures for various ferromagnetic and non-ferromagnetic metal layers such as Co (Cobalt), Ni (Nickel), Mo (Molybdenum), Ru (Ruthe-

nium), Os (Osmuim), Ir (Iron), Cu (Copper), Ag (Silver), etc. The magnetoresistance value depends on the thickness of these metal layers and the type of pair of ferromagnetic and non-ferromagnetic materials. In all those materials, it was found that the GMR is associated with changes in the relative orientation of the magnetization in the successive ferromagnetic layers. This effect is termed as spin-value magnetoresistance [85]. This effect can be explained using a triple-layered metal structure which has two identical ferromagnetic layers sandwiching a non-ferromagnetic metal layer. When the two ferromagnetic layers are magnetized parallel to each other, a particular spin type (either spin-up or spin-down) of electrons can easily move through the non-ferromagnetic layer in an ordered way, exhibit a low resistance. On the other hand, in the case where those two ferromagnetic layers are anti-parallel to each other, both spin-up and spin-down electrons undergo collisions in either of the ferromagnetic layer (depending on the type of spin), resulting in a higher resistance. Hence, the parallel and anti-parallel magnetic configuration *opens* and *closes* the flow of electrons, respectively, acting as a valve. There is another effect, in which the non-ferromagnetic layer of the spin valve structure is replaced with an insulating layer known as *tunnel effect* [85]. In this structure, the electron move from one ferromagnetic layer to the other by a tunnel effect while conserving their spin. This effect was first proposed by M. Julliere in 1975 [87] and experiments were conducted on Fe/Ge/Co multilayer at 4.2K. Nevertheless, it did not attract much attention because the magnetoreistance value obtained was very small around 14 % [87]. This effect draws a significant attention after the discovery of GMR, that paved the way towards the development of Magnetic Tunnel Junction, which is the key element of the *Magnetic Random Access Memories* (MRAM) [41].

The discovery of GMR had triggered an immense research activities as this has the potential to be used for many useful applications. For instance, this effect was immediately adopted in hard disk driver industries for the read sensors as it allows more data to be packed on computer disks [41, 88–90]. The magnetoresistance head structure employed by IBM for its hard disk structure [91]. In general, a disk consists of tiny regions of magnetization that is covered with a thin film of magnetic materials and the data is stored (in term of zeros and ones) in the form of the direction of the magnetization of these regions. Alternatively, the data is read by sensing the magnetic fields using the head that flies at a constant height above the magnetic domains. With the increase in density, these regions of magnetization get smaller, the sensing mechanism to read the data becomes challenging. The GMR effect significantly increases the read margin (i.e. change in resistance between the two states), that improves the sensing methodology and hence, allows more data to be stored. Additionally, the discovery of GMR effect is also important for various other applications such as in the development of MRAM, magnetic-filed sensors and isolators [92, 93].

### 2.2.2 Magnetic Tunnel Junction

The next big step for the spintronic technology was the development of *Magnetic Tunnel Junction* (MTJ), which is the storing component in MRAM. In MTJs, the values are stored in terms of resistance states. An MTJ cell consists of the two ferromagnetic layers separated by an insulator layer. The typical structure of the MTJ cell is shown in Figure 2.5. In general, this insulator layer, which is also referred to as the barrier layer or tunnel barrier, is a very thin layer (a few nanometer) so that electron can easily tunnel through this layer, resulting in electrical conduction. This conduction using tunnelling process depends the relative magnetic orientation of the two ferromagnetic layers. When the magnetic orientation of these two layers are *Parallel* (P) to each other, it exhibits a low resistance value. On the other hand, the resistance value is high in the case when the magnetic orientation of these two layers are *Anti-Parallel* (AP) to each other. Considering this fact, the conductance becomes as described in [94]:

Figure 2.5: A typical MTJ cell structure

$$G(\theta) = \frac{1}{2}(G_P + G_{AP}) + \frac{1}{2}(G_P - G_{AP}) \cdot cos\theta \qquad (2.1)$$

where $\theta$ represents the relative magnetic orientation of the two ferromagnetic layers as described by Figure 2.5, $G_P$ and $G_{AP}$ are the conduction when the magnetic orientations for the two ferromagnetic layers are parallel (i.e. $\theta = 0°$) and anti-parallel (i.e. $\theta = 180°$), respectively. The magnetoresistance observed due to the tunnelling process (known as *Tunnelling Magneto-Resistance* or simply TMR) for an MTJ cell is defined as [94]

$$TMR = \frac{G_P - G_{AP}}{G_{AP}} = \frac{R_{AP} - R_P}{R_P} \qquad (2.2)$$

where $R_P$ and $R_{AP}$ are the resistance of the parallel and anti-parallel magnetization states, respectively.

The origin of the TMR effect can be understood with a schematic diagram as depicted in Figure 2.6. During tunnelling, the electron spin orientation is preserved and only the electron with the same spin orientation as the spin-band can tunnel to the next electrode. Here Figure 2.6(a) shows the parallel state magnetization, which means the majority states can fill the



(a) Parallel magnetization          (b) Anti-parallel magnetization

Figure 2.6: Illustration of the TMR effect for an MTJ cell. The electron with the same spin as the sub-band can only tunnel through. The impact of tunnelling electrons during: (a) parallel magnetization, and (b) anti-parallel magnetization of the the two ferromagnetic layer has. Here dashed line show spin conserved tunnelling.

Figure 2.7: Trend of TMR using amorphous aluminium oxide and crystalline magnesium oxide as tunnelling barrier (data obtained from [8–13] and [14–20]).

majority states of the next electrode and vice-versa. In this figure, the majority and minority tunnelling through electrons are shown with thicker and thinner arrows, respectively. The parallel magnetization results in a good band matching, which means large conductance and exhibit a small resistance. On the other hand, for anti-parallel magnetization as shown by Figure 2.6(b), the electron with minority spins fills the majority sub-band, while minority spins fill the majority sub-bands in the next electrode. The anti-parallel magnetization results in a poor band matching, which means attenuated current and exhibit a high resistance. The TMR effect was also expressed by the Julliere's model [87]:

$$TMR = \frac{2P_1P_2}{1 - P_1P_2} \qquad (2.3)$$

where $P_1$ and $P_2$ are the polarisation factor of the two ferromagnetic electrode. The polarization factors can be defined as [95]:

$$P = \frac{N_{majority} - N_{minority}}{N_{majority} + N_{minority}} \qquad (2.4)$$

where N is the spin-resolved density of the respective states for the two ferromagnetic electrode.

The TMR effect using an amorphous aluminium oxide ($AlO_x$) as tunnelling barrier was first demonstrated by J.S. Moodera [96] and T. Miyanzaki [97] independently in 1995. Further, J.S. Moodera and T. Miyanzaki conducted experiments using MTJ stacks of CoFe/$Al_2O_3$/Co or CoFe/$Al_2O_3$/NiFe and Fe/$Al_2O_3$/Fe, respectively. The TMR value observed during that time using these experiments was slightly more than 10 %, which latter improved upto 70 % using CoFeB/$Al_2O_3$/CoFeB MTJ in the year 2004 [11]. Recently, Kap Soo Yoon demonstrated a high magnetoresistance of about 220 % using CoFeN/$AlO_x$/CoFeN MTJ [8]. It was observed that the TMR can not be improved using amorphous aluminium oxide as tunnelling barrier after a certain level because of its unsymmetrical atom organisation, that leads to an incoherent tunnelling. The TMR value is considerably increased with the introduction of the crystalline magnesium oxide (MgO) as tunnelling barrier, which delivers the structurally ordered junctions [14–20]. The trend of the growth of TMR (in % value) for both $AlO_x$ and MgO is illustrated in Figure 2.7. This begins with a prediction made by W.H. Butler in year 2001 [98].

The significance of the TMR value is for the read operation, the higher the TMR value more efficient the read is. The stored value in MTJs can be read by passing a current through the stack. The stored value corresponds to the resistance state of the MTJ and the resistance state

is decided based on the magnetic states of the ferromagnetic layer. The current value is then sensed using a sense amplifier, that generates the output value after comparing the current with that of a reference value. In order to demonstrate this device as a true storing cell, its magnetization should be controllable. The mechanism to control the ferromagnetic layer is known as *storage* (or write operation or switching) operation. To make MTJ for storage, one layer has to be freely rotated and the other can act as reference layer whose magnetization is always fixed.

## 2.2.3 Magnetic switching

In order to make MTJ a memory cell, another important characteristics of MTJ besides TMR is to write (or store) a given value. For that purpose, a magnetic switching mechanism needs to be developed for an MTJ cell. Many switching mechanisms have been developed in the past two decades, those we will discuss next.

### Field Induced Magnetic Switching

The *Field Induced Magnetic Switching* (FIMS) approach was used for the conventional MRAM technology. In this approach, the magnetic orientation of the free layer of the MTJ is switched using a two orthogonal magnetic fields. A cross point architecture of a FIMS based MRAM is illustrated in Figure 2.8. It consists of an array of MTJs sandwitched between the two conducting lines, namely, bit-line and digit-line as shown in the figure. In order to write a value or to aligh the magnetization of the free lines in a required orientation, current pulses are sent through bit-lines and data-lines. These current flow generate magnetic fields along the conducting lines and an MTJ cell at the intersection of these two fields is selected for the write operation. This switching mechanism is based on a coherent reversal magnetic orientation proposed by Stoner-Wohlfarth model. The magnitude of the magnetic fields required for the magnetic switching is described by an asteroid equation:

$$H_h^{\frac{2}{3}} + H_e^{\frac{2}{3}} \geq \left[\frac{2K}{M_s}\right]^{\frac{2}{3}} \tag{2.5}$$

where $H_h$ and $H_e$ are the two magnetic fields, $K$ is is an effective anisotropy constant and $M_s$ is the saturation magnetization. Using this approach, the magnetic switching process does not involve any actual movement of electrons, resulting in almost no wear-out. Another advantage is that this device is non-volatile, meaning that the magnetic polarization does not leak, so the information remains stored even the power is turned off. However, one of the biggest drawback of this technology is that there are many half-selected bits in each write operations and the possibility of unnecessary switching of such bits is very high. Moreover, due to a narrow operating window, it is challenging to tune the fields in order to switch the required number of bits.



Figure 2.8: Illustration of a cross-point architecture using field induced magnetic switching based devices.

To overcome these issue, Savtchenko proposed a toggle switching approach [99, 100] in which the memory cell is typically oriented at 45° with respect to the the two conducting line. A synthetic ferromagnetic layers of MTJ cell is considered which is separated through a thin Ru spacer layer. In toggle MRAM, an offset in time has to be established between the two current pulses such that they produce a 180° rotation in magnetization of the free layer of the selected MTJ cell. A single current field leads to a small magnetic rotation of the free layer, and it is fully reversible when another current pulse is applied. This way, half-selected bits issue is completed eliminated. Nevertheless it requires a read before every write to know the content of the bit, because this approach toggles the bit to the opposite magnetic state regardless of existing state. Additionally, FIMS has a high energy consumption as it requires a high current flow through the conductive lines to generate magnetic fields.

### Thermally assisted Magnetic Switching

In order to further improve the magnetic switching process, a *Thermally Assisted Switching* (TAS) approach was proposed [21, 101–103]. In this approach, a heating mechanism is utilized that not only accelerates the magnetic switching process, but also eliminates the requirement of one of the write fields. Here, the MTJ stack is slightly modified by including an additional layer of biased exchange composed of an anti-ferromagnetic material to the storage layer. The switching process of the TAS approach as proposed by [21] has total of three steps, as depicted in Figure 2.9. First is heating, that is generated by passing current in the write lines through the MTJ cell as shown in the figure. Once heating exceeds the blocking temperature of the storage layer, then the writing process is initiated by applied external magnetic field $H_{SW}$. This external magnetic field is generated by passing current through the write field as in FIMS approach. Finally, the cooling step in which the device starts cooling down which is performed in the presence of the magnetic field. This approach has several advantages such as better bit selectivity during write operation, power consumption is improved as it works with a single external magnetic field and data retention ability is improved due to exchange bias of the storage layer. However, its total switching time is too long and consumes very high energy, although lower than the conventional MRAM, because it requires an additional current to generate heat before the actual write process starts.

### Spin Transfer Torque Magnetic Switching

Previous magnetic switching approaches consume very high energy due to the requirement of the immense current to generate external magnetic fields and/or heat. Moreover, they require a long delay to complete the entire switching process. With such a high access energy and latency values, it is nearly impossible to employ these devices for a low-power and high-performance



Figure 2.9: Demonstration of switching steps in a thermally assisted magnetic switching device [21].

(a) Interaction of electrons with the magnetization of free layer

(b) Switching principle

Figure 2.10: Spin transfer torque switching mechanism.

applications. To address these issues, a new switching mechanism was discovered in which spin polarized current is passed through the MTJ stacks to perform the magnetic switching, known as *Spin Transfer Torque* (STT) effect. In this effect, the spin polarized current carries angular momentum, that can be transferred to the magnetization. The STT effect was first experimentally demonstrated on Co/Cu/Co stack in the year 2000, which is later also proven working on MgO based tunnel barrier.

The basic principle of the current-induced torque is illustrated in Figure 2.10(a). This figure shows the interaction of the spin of the electron with the magnetization of a ferromagnetic layer, that results in a diversion in the electron spin after passing through the ferromagnetic layer. This change in the direction of the spin angular momentum of the electron leads to a torque on the magnetization of the ferromagnetic layer, which eventually resulting in a magnetic switch. The STT switching process is shown in Figure 2.10(b). In order to perform 'AP'→'P' magnetic switching in the free layer, electron should flow from the reference layer to the free layer as depicted in the figure. The electrons those have the same spin directions as that of the magnetic orientation of the reference layer, tunnel through barrier layer and enter into the free layer. As a consequence, the spin angular momentum of the electron is transferred to the magnetization of the free layer, which results into a torque. If the generated torque is sufficiently strong exceeds the threshold value (known as the *critical current* for 'AP'→'P' switching), the magnetization of the free layer is switched. On the other hand, the electrons should flow in the reverse direction (i.e. from free layer to the reference layer) to perform 'P'→'AP' magnetic switching in the free layer. As shown in the figure, after passing the free layer, only those electrons which have the same spin as that of the magnetization of the reference layer flow through. Whereas ones with opposite magnetic directions are reflected back into the reference layer from the boundary between tunnel barrier and the reference layer. The spin angular momentum of these electrons are transferred to the magnetization of the free layer, that generates a torque and resulting in a magnetic switch if the torque exceeds the threshold value (known as the *critical current* for 'P'→'AP' switching). The critical current density of the two switching can be expressed as

$$J_{co}^{P \to AP} = \alpha \gamma e M_s t / (\mu_B g(0)) \cdot [(H_{ex} + H_{dip}) + (H_{ki} + H_d)] \tag{2.6}$$

$$J_{co}^{AP \to P} = \alpha \gamma e M_s t / (\mu_B g(\pi)) \cdot [(H_{ex} + H_{dip}) - (H_{ki} + H_d)] \tag{2.7}$$

where $\alpha$ is the Gilbert damping constant, $\gamma$ is the gyromagnetic factor, $e$ is the electron charge, $t$ is the thickness of the free layer, $\mu_B$ is the Bohr magneton constant, $M_s$ is the saturation

magnetization. The $H_{ex}$, $H_{ki}$, $H_{dip}$ and $H_d$ are the in-plane applied, in-plane anisotropy, dipole fields from the pinned layer acting on the free layer and effective demagnetization field, respectively. The $g$ is the STT efficiency that can be expressed for an angle between the magnetization of the free layer and the reference layer ($\theta$) as:

$$G(\theta) = P/\left[2(1+P^2)cos\theta)\right] \tag{2.8}$$

The effective demagnetization field can be expressed as

$$H_d = M_s/\mu_0 - H_{kp} \tag{2.9}$$

where $H_{kp}$ is the perpendicular anisotropy field.

This approach has several advantages compared to previous magnetic switching schemes. For instance, this scheme does not require any magnetic field to switch the magnetization of the free layer, which also significantly improves overall energy consumption and area. Moreover, the selectivity of bits is not at all a problem in STT based switching approach as current has to pass through the stack to perform switching. The switching current density can be reduced with the reduction of the size of the MTJ, hence has better scalability. However, in this technology the switching current is still too high, that needs to pass through the MTJ stack for a long duration to make a magnetic switch. This leads to high energy consumption, high switching delay and several reliability issues such as *Time dependent Dielectric Breakdown* (TDDB). These issues are resolved by a large extent by using a out-of-plane equilibrium orientation of the magnetization of the free layer of the MTJ cell, that is termed as *perpendicular STT* switching. The reason of the requirement of high current and long duration for the in-plane devices is that they have to overcome a very large out-of-plane demagnetization field during switching process. Whereas, for *perpendicular STT* (STT-P), since the magnetic switching is perpendicular, the demagnetization field factor is naturally eliminated. Due to this, it requires less torque for magnetic switching. Thus switching can be possible with considerable low current and in less time, resulting in better energy compared to in-plane STT switching. The critical current density of the two switching for STT-P can be expressed as

$$J_{co}^{P \to AP} = \alpha\gamma eM_st/(\mu_Bg(0)) \cdot [-(H_{ex}+H_{dip})+(H_{kp}-M_s/\mu_0)] \tag{2.10}$$

$$J_{co}^{AP \to P} = \alpha\gamma eM_st/(\mu_Bg(\pi)) \cdot [-(H_{ex}+H_{dip})-(H_{kp}-M_s/\mu_0)] \tag{2.11}$$

**Spin Orbit Torque Magnetic Switching**

Alternative to STT, recently a new current induced approach is discovered, named as *Spin Orbit Torque* (SOT). In this technology, the magnetic torques are generated based on the spin-orbit interaction that rotates the magnetic orientation of the free layer of the MTJ. In fact, the magnetic stacks are placed on the top of a heavy metal such as platinum (Pt) or tantalum (Ta), which are capable of generating strong spin-orbit interaction. The *Spin Hall effect* (SHE) or/and the *Rashba effect* are responsible for the current induced magnetic switching. The SHE is a electron transport phenomenon in which spin accumulation appear on the lateral surfaces of the conductive metal [104–106]. Due to this effect electrons with opposite spin drift in opposite direction as shown in Figure 2.11(a). This phenomenon is also applied for inverse-SHE as illustrated in Figure 2.11(b).

Like in STT, the switching current does not flow through the MTJ stack in SOT based MTJ, instead it passes through a heavy metal layer to generate spin-orbit interactions. Therefore, it requires an additional terminal to provide a bidirectional current path. A three terminal SOT based MTJ cell is illustrated in Figure 2.11(c). As shown in the figure, when current is

(a) Demonstration of Spin Hall Effect

(b) Demonstration of Inverse Spin Hall Effect

(c) Spin orbit torque switching principle

Figure 2.11: Spin orbit torque switching mechanism

passing through the heavy metal layer, spin-orbit coupling causes spin-dependent scattering or deflection, resulting in a spin angular momentum along the direction perpendicular to the electron flow. This leads to a spin accumulation at the boundaries of the metal layer which can be extracted by the adjacent ferromagnetic layer (free layer of the MTJ) and exerting a torque, results in its magnetization switch. In SOT technology, both in-plane and perpendicular magnetic switchings are possible. For the in-plane magnetic switching, it has to overcome the demagnetization field similar to that of the STT. On the other hand, for perpendicular switching, an additional magnetic field is required that is placed at the top of the MTJ stack in order to achieve a deterministic switching.

### 2.2.4 MRAM characteristics

#### Asymmetric switching

In STT-MRAM, write current flows through the bit-cell for a certain duration to flip its magnetization. This current has a bidirectional path and the direction depends on the input value to be written in the bit-cell. Moreover, it has an asymmetry in the switching latencies for writing '1' and writing '0' values [45–48]. Due to this asymmetrical behavior, writing 'P' configurations (value '1') is fast which we refer to as fast-writes and writing 'AP' configurations (value '0') is slow (slow-writes). This asymmetry in the bit-cell is categorized into two parts: the asymmetry due to the MTJ cell and the asymmetry due to the bit-cell access transistor.

The MTJ cell has an inherent torque asymmetry because of its two different mechanisms to flip the magnetic orientation. The transition to the 'P' state is based on the same spin direction electrons as the magnetic orientation of the RL, whereas the switching into the 'AP' state is due to the opposite spin direction electrons that are reflected back at the boundary of the oxide layer and the RL [107]. This asymmetry of the MTJ switching behavior depends on the spin-transfer efficiency factor of the ferromagnetic layers [45]. For the MTJ model as specified in [108], write latency for writing 'P' and 'AP' configuration is 3.9 ns and 5.7 ns, respectively, as demonstrated in Figure 2.12.

The other source of asymmetry is the voltage degradation of the access transistor (equal to the threshold voltage of NMOS). This voltage degradation reduces the write current for the slow-write operations (i.e. 'AP'), and consequently the switching delay for that kind of write operations increases significantly. As illustrated in Figure 2.12, the write latency for an 'AP' configuration increases from 5.7 ns to 9.7 ns after adding the access transistor, while the latency for the fast-write remains the same.

Figure 2.12: Demonstration of asymmetric switching behavior of the MTJ using waveforms.

## Stochastic switching

The switching of the MTJ cell is stochastic in nature due to the random thermal fluctuations. This means that the switching delay of the magnetization is not deterministic. However, the write period has a fixed duration, and thus incomplete transitions due to the stochastic write behavior cause write errors. The *Write Error Rate* (WER) of a bit-cell or the overall write duration of $t$ is expressed as [49]

$$P_{WER}(t) = 1 - exp \left[ \frac{-\pi^2 \Delta (\frac{I}{I_c} - 1)/4}{\frac{I}{I_c} e^{2\alpha\gamma H_K t (\frac{I}{I_c} - 1)/(1+\alpha^2)} - 1} \right] \qquad (2.12)$$

In this equation, $I$ is the write current, $I_c$ is the critical current, $\alpha$ is the Gilbert damping constant, $\gamma$ is the gyromagnetic factor, $\Delta$ is the thermal stability factor and $H_K$ is the effective anisotropy field. In general, the typical value for the targeted WER for the STT-MRAM is $10^{-9}$ with *Error Correction Code* (ECC) and $10^{-18}$ without ECC correction [109]. It is inferred from the above equation that the switching probability of the cell increases by extending the duration of the write period [110, 111]. For instance, using MTJ model as specified in [108], WER of $10^{-9}$ can be achieved with a write period of 18 ns. Please note that the write period is decided based on 'P'→'AP' switching delay as it is the worst of the two switching delays. Using the above equation, the write latency distribution for both 'AP' and 'P' switching delays for a single bit-cell are shown in Figure 2.13. As it can be seen in the figure, the 'AP' switching delay has a wider distribution compared to the 'P' switching delay. Furthermore, these distributions have a very long tail, especially for the 'AP' switching, resulting in the aforementioned long write period (see Figure 2.14).



Figure 2.13: Switching distribution of an MTJ for a single bit-cell [22].

Figure 2.14: Write energy and write error rate for various write periods for a single bit-cell [22]. This illustration is for the 'AP' switching, which represents the worst case switching behavior.

## Read Disturb

In STT-MRAM, the content of a bit-cell can accidentally change during a read operation which is known as read disturb. This is due to the fact that the read current shares one of the write current paths in STT-MRAM. However, the read current is around 5-10 times lower than the critical write current (minimum current required to flip the bit-cell at a certain write period and write error rate). Nevertheless, this low read current induces a magnetic disturbance in the MTJ cell which may lead to a flip of the magnetic orientation. Since the read current is unidirectional, the flip can only happen in one direction, i.e. either from 'AP'→'P' or the other way around. As a result, also the resistance changes, which in turn affects the read current.

The switching probability due to read disturb is given by the following equation [109]:

$$F_{rd} = 1 - e^{-\dfrac{t_{read}}{\tau_1 e^{\Delta(1-I_{read}/I_{C0})}}} \tag{2.13}$$

where $\Delta$ is the thermal stability factor, $I_{read}$ is the read current, $I_{C0}$ is the write critical current, and $t_{read}$ is the read period. For STT-MRAM, the typical read disturb probability for a single read event is in the range of $10^{-23}$ to $10^{-21}$ [109].



Figure 2.15: Conceptual diagram of read disturb while reading 'AP' state (solid line indicates actual current flow)

### Degradation effect

The MTJ cell consists of very thin layer of MgO tunnel barrier (thickness of around 1.2 nm), which can lead to a breakdown under high stress conditions known as *Time Dependent Dielectric Breakdown* (TDDB) [112, 113]. This is a severe reliability challenge for the MTJ cells and can considerably limit their lifetime. The time to breakdown value depends on the amount of current flowing through the device and its duration. It can be modeled in a similar way as the gate dielectric breakdown modelling in MOSFET, which is given by following equation [114]:

$$W_{hard\_breakdown} = \beta ln\ t - \beta ln\ \alpha \tag{2.14}$$

where $\alpha$ and $\beta$ are current-dependent parameters, and $t$ is the duration of the current flow.

### Process variation

The development of spintronic based memory requires two different fabrication processes, i.e., a *magnetic process* and a *CMOS process*. Any variations in these processes affect the cell properties. Due to the manufacturing process, the MTJ cell has mainly variations in area, TMR and oxide thickness. This results in disturbances in the magnetization of the cell, affecting the write current that leads to a deviation of the switching delays from its mean [115]. The additional timing margin for the correct functionality of the memory in the presence of process variation could be as large as 2 X [116]. Therefore, it is critical to take the impact of process variation into account.

### Defects in MTJs and fault modeling

The MTJ device uses different materials and processes for manufacturing compared to CMOS. Due to the complexity of these fabrication processes and the interdependency of magnetic materials, MTJ cells are subject to various and new failure mechanisms [25, 114, 117–119]. For instance, during the ion beam etching process, due to sputtering effects, the sputtered atoms re-deposit at the MTJ sidewall leading to a barrier short [118]. As a consequence, the isolation of the ferromagnetic layers is damaged, resulting in a very low resistance value [114]. In this case, although current can flow through the device, the cell no longer behaves as an MTJ cell. On the other hand, an MTJ cell can also have an open connection due to internal damage, which delivers a very high resistance value. As a result, current cannot flow through the device because of a discontinuity in the design.

There are some cases, where the MTJ cell resistance values are not so severely affected as in open and short faults, however, their values can easily influence the sense amplifier so that an incorrect output can be generated. For instance, during fabrication, the magnetization of



Figure 2.16: Classification of MTJ fault models.

the FL can be permanently fixed to either 'P' or 'AP' configuration [120]. Another possibility is that the switching margin and/or the current value are not sufficient enough to flip the magnetization of an MTJ cell [25]. This is due to defects or the impact of process variation. In addition to manufacturing defects, MTJs are also vulnerable to runtime failures such as read disturb, retention failures and back-hopping [119, 121, 122].

We can broadly classify all these MTJ defects into the following four fault models (as illustrated in Figure 2.16):

- *Short fault*: The two ferromagnetic layers, FL and RL are connected.

- *Open fault*: Discontinuity in the device.

- *Stuck-at-P fault*: MTJ magnetization is permanently or temporally locked to the 'P' state.

- *Stuck-at-AP fault*: MTJ magnetization is permanently or temporally locked to the 'AP' state.

## 2.3 Memory architectures using spintronic technology

Like any other memory technologies, MRAM bit-cell also comprises of a storage element and its access mechanism. In MRAM technology, the storing element is an MTJ cell that stores the value in term of resistance states as described previously. In order to access certain MTJ cells for the random access memory operations (such as store or read operations), one or more CMOS-based transistors are required. This is due to the fact that the MTJ cells are highly compatible with the CMOS technology because of their suitable resistance values. Next, the bit-cell architectures of STT-MRAM and SOT-MRAM technologies are explained.

**Bit-cell architecture**

The fundamental building block of a MRAM is the MTJ based memory cell. Figure 2.17(a) shows a typical STT-MRAM bit-cell architecture. It consists of an MTJ cell and an access transistor (refereed as 1T1M architecture). It has total three terminals namely, *Source-line*, *Bit-line* and *Word-line*. The access transistor can be connected to the Free Layer or the Pinned layer of the MTJ cell. More commonly, it is connected to the Pinned layer, therefore known as *standard connection*, on the contradictory if it is connected to the Free layer, it is known as



(a) STT-MRAM bit-cell      (b) SOT-MRAM bit-cell

Figure 2.17: Typical bit-cell architectures for STT-MRAM and SOT-MRAM memories.

*reverse connection.* The other end of the MTJ cell and the access transistor are connected to the Source-line and Bit-line terminals, respectively. These two terminals facilitates the read-write current flows during the corresponding memory operations. On the other hand, the gate of the access terminal is connected to the word-line terminal, which is activated based on the output of the row decoder.

The bit-cell architecture for SOT-MRAM is demonstrated in Figure 2.17(b). As shown in the figure, it consists of an SOT based MTJ cell and two NMOS-based access transistors. This architecture has total of four terminals namely, *Bit-line*, *Source-line*, *Read Word-line* and *Write Word-line*. One access transistor ('N1') is connected to the read terminal of the SOT based MTJ cell, which is activated only during the read operations. Whereas, the other access transistor ('N2') is connected to the write terminal of the MTJ cell, which is activated only during the write operations. The other end of both transistors are connected to the bit-line terminal. At a time only one transistor is activated based on the output of the row decoder depending on the memory operation. For instance, for read operations the access transistor 'N1' is activated so that read current can flows through the bit-line to the source-line terminals. On the other hand, for write operations, the other transistor 'N2' is activated so that write current can flow between the source-line and bit-line terminals. Due to device-level phenomenon, for a given memory operation, a certain amount of unwanted current (known as *sneaky current*) also flows out through one of the other terminal which is not in use. For instance, during write operation, a certain amount of current can flow through the MTJ stack if the access transistor is not connected at the read terminal of the MTJ cell. Such a sneaky current can flow in both read and write operations, that can not only disturb the operating cell but also erode the content of the neighbouring cells.

### Array architecture

The array organization of any memory has primarily two requirements: (1) to provide accessibility to bit-cells, and (2) to facilitate the memory operation. The bit-cells are accessed using a row decoder, whereas memory operations are performed using their respective read-write circuitries. the basic array organization for both STT-MRAM and SOT-MRAM is shown in



(a) Sub-array architecture using STT

(b) Sub-array architecture using SOT

Figure 2.18: Typical sub-array architectures for STT-MRAM and SOT-MRAM memories.

Figure 2.18. In general, memory bit-cells are accessed using the decoder circuit based on the applied memory address. These memory addresses are first latched and activate the corresponding word-line using decoder logics, and a driver circuitry is used to drive this signal. A *pre-decoding scheme* is employed in these decoding logics in order to improve the area. The decoder circuits have extensively explored in the conventional memory technologies and similar circuits can also be applicable to MRAM designs. In STT-MRAM, decoder logics directly activate the word-line signal as illustrated in Figure 2.18(a). On the other hand, SOT-MRAM bit-cells have separate word-lines for read and write accesses (i.e. $WL_R$ and $WL_W$), which require a mux to activate the respective word-line through the decoder logic as shown in Figure 2.18(b). The control-signal of the mux is controlled using the *write enable signal*, which is the memory read-write operation control signal. When the write enable signal is high means write operation is performed and vice-versa. The read-write circuities are also activated based on the control signal. The MRAM specific read-write circuitries are discussed next.

**Write circuit**

For the write operation, a bidirectional current path has to be established for every bits. The schematic diagram of a typical write circuitry is shown in Figure 2.19(a). As shown in the figure, write circuit has total four terminal namely, write_enable, data_in, and the two output terminals. The two output terminals need to be connected to the bit-line and source-line terminals of the bit-cell. The write_enable signal activates the write operations, i.e. it is 1 for the write operation otherwise it remains 0 for read or no memory operations. The data_in is another input signal, from where input data has to be provided, which eventually decides the direction of the current flow. This is because, the current direction makes the magnetization switching in the required orientation (either 'P' or 'AP'), that represents different resistance states, hence distinguishable two logic values. (about write buffers, may be in array modelling part).

**Read circuit**

For the read operation, a small unidirectional current has to pass through the MTJ stack. This current can be sensed using a sense amplifier in order to distinguish between the two resistance



(a) Write

(b) Read

Figure 2.19: Read and write circuitries for both STT-MRAM and SOT-MRAM.

state after comparing it with the reference signal. A pre-charge based current sense amplifier is shown in Figure 2.19(b). It mainly consists of an equalizer circuit, a reference bit-cell and a sensing circuit. In the equalizer circuit, the two output nodes (q1 and q2) need to be at the same potential before the sensing operation begins. This is achieved using an equalizer circuit which is controlled using a pre-charge (PC) signal. The equalizer circuit is active when the PC signal is '0' and becomes inactive when PC is '1' during the actual sensing process. The reference bit-cell circuit consists of a set of four MTJ and a NMOS transistor as shown in the figure. The four MTJs connected in such a way that the effective resistance value is the middle of the two resistance states that a single MTJ can take i.e. $R_P + R_{AP}/2$ where $R_P$ and $R_{AP}$ are the resistance values during 'P' and 'AP' magnetization state. Moreover, the sensing circuit consists of two back-to-back connected inverters. The sensing circuit is stimulated with the deactivation of the pre-change signal and activation of the word-line signal. The reference word-line signal is activated using *address transition detection* circuit. During the read evaluation, the two output nodes read and $\overline{read}$ try to become stable based on current values in the two branches, that in turn depends on their resistance values.

## 2.4 Normally-off computing enabled by non-volatile spintronic based latches and flip-flops

Flip-flops and latches are the basic building blocks of digital electronic systems. These are the storing component in a sequential logic, that are extensively used in a system-on-chip. In such device, data can be stored as a state in a sequential logic, value of a counter or a memory element to store any other piece of information. Additionally, flip-flops are adopted in a synchronous computation system as a cycle boundary in order to evaluate their path delays. In general, flip-flops are composed of a pair of latches, whose functionality is controlled with the different levels of a clock signal. A typical flip-flop architecture is shown in Figure 2.20. As shown, the storing mechanism in a latch is organized using a two back-to-back connected inverters similar to an SRAM bit-cell structure. These two inverters always have to be connected to a supply voltage to retain its stored data, even when the design is operated in a standby mode. Since these flip-flops can not be power-gated during standby mode, it accelerates the leakage consumption of the design. In order to deal with this issue, spintronics storing devices such as STT and SOT are very effective as these are non-volatile and can be easily power gated.

Many non-volatile shadow flip-flop architectures have recently been introduces which exploit the normally-off and instant-on attributes of this technology. The block diagram of a typical shadow flip-flop architecture is shown in Figure 2.21. It consists of three components, namely, *master latch*, *slave latch*, and *non-volatile shadow latch*. As shown in the figure, the master and slave latches are the same as in conventional CMOS flip-flop design and the non-volatile



Figure 2.20: A typical CMOS-based flip-flop architecture

(a) A typical shadow flip-flop architecture

(b) Shadow latch using STT-based MTJs

(c) Shadow latch using SOT-based MTJs

Figure 2.21: Illustration of a non-volatile shadow flip-flop architecture using MTJs.

latch component is connected to the slave-latch. This type of design have an additional pin named *power-down* (PD) pin. During standby mode, with the activation of this pin, the slave data value is stored in the shadow latch. After that, the entire design blocks including flip-flop designs can be power gated to reduce the static power consumption. On the other hand, during wakeup, the non-volatile shadow latch content is read and re-stored into the slave latch, so that the normal operation can be resumed.

The block diagrams of a typical non-volatile latch architecture using STT and SOT based storing MTJ cells are shown in Figure 2.21(b) and Figure 2.21(c), respectively. It consists of two MTJ cells and CMOS-based read and write components as shown in figures. The write or store operation is activated based on 'PD_wr' signal and the value to be stored is decided based on the direction of the current flow which is controlled using 'D' signal. The two MTJs should always store the opposite magnetizations which assists the read process by providing a self-referencing structure while sensing the resistance differences. On the other hand, the read is activated using 'PD_rd' signal and the value is read using current sensing mechanism, which is available at mtj_read and $\overline{mtj\_read}$. The 'PD_rd' and 'PD_wr' signals are generated using the 'PD' pin.

## 2.5 Memory architecture-level evaluation

Before going into details of memory architecture evaluation, a qualitative analysis of MRAM technology is performed in comparison to other conventional and emerging technologies. Figure 2.22(a) describes the comparison of MRAM technology with respect to the convention memory technologies. As shown in the figure, MRAM is highly scalable compared to all other conventional memories. Moreover, the performance of MRAM and endurance[2] level is similar to that of SRAM. Nevertheless, the write access energy consumption is relatively higher compared to that of SRAM, which is at the same level as DRAM and much better than Flash. This is because, DRAM requires a periodic refresh and Flash write operations require the charge pump to supply a high voltage (5V to 20V) [123]. The MRAM technology can be as high density as DRAM as shown in the figure. The NAND-Flash has highest density among all because of its multi-level cell capabilities. Moreover, MRAM can have the similar level of data retention[3] as Flash because its bit-cells are highly non-volatile in nature. As per ITRS, MRAM can retain their data for a duration of more than 10 years. However, in reality for high

---

[2]Endurance is the number of cycles that a storing cell can sustain.

[3]Retention time refers to the data retaining capability of the storing cell, regardless of their powered-on or powered-off.

(a) Conventional memories  (b) Emerging Memories

Figure 2.22: Comparison of MRAM technology with conventional and emerging memories

performance memories such as caches in standard computing system, such a high retention duration is not required. Therefore, to target such memory systems, the retention values of MRAM can be relaxed, that can result in improvement in the write access latency and energy. Additionally, MRAMs are immune to the soft-errors due to radiations unlike to conventional memories.

On the other hand, the comparison of MRAM with other emerging memories are shown in Figure 2.22(b). All these memory technologies are proven scalable, but MRAM is the only memory technology among them, whose endurance level is equal to that of SRAM. Additionally, MRAM and RRAM have almost similar access latency and write energy values, which are relatively better than those of PCRAM and FRAM. The array density of the PCRAM and RRAM is better than MRAM technology. However, PCRAM and RRAM have another reliability challenge, that is, the functionality of their cells are material dependent whose properties are significantly changing with the temperature. On the other hand, FRAM has a destructive read cycles, means a write back has to be performed after every read operations. The MRAM technology does not have such limitations.

In summary, the MRAM is a highly scalable technology that can be as fast as SRAM, as dense as DRAM and as volatile as Flash [23, 124–128]. Moreover, it is more reliable and has the best latency and energy efficiency compared to other emerging memory technologies.

The memory architecture-level evaluation is performed using NVSim [129]. NVSim is a circuit-level module for *Non-Volatile Memory* (NVM) performance, energy and area estimation which uses an empirical modelling methodology based on the well-known CACTI [130]. It estimates the access time, access energy and the total area of NVM chips with different options before the fabricating of actual chip. The bit-cell level information and the desired memory architecture information like capacity, data width, associativity, local and global wire type, routing type etc., are inputs to this tool. It also supports memory array organizations which have three hierarchy levels i.e. bank, mat and subarray. Furthermore, NVSim has optimization settings such as buffer design optimization, optimization target, and various other design constraints. The sizing of the periphery circuit is done by optimized models i.e. latency, area and balance optimization. NVSim generates results such as write access latency for both set and reset, read access latency, read dynamic energy, write dynamic energy for both set and reset, leakage power and area for the given memory configuration. Using this platform, a quantitative

(a) Area

(b) Latency

(c) Energy

(d) Leakage

Figure 2.23: Memory scaling results for SRAM, STT-MRAM and SOT-MRAM [23].

comparison for SRAM, STT-MRAM and SOT-MRAM is performed, and results for the same are as depicted in Figure 2.23.

### System-level evaluation

The output data from NVSim (as described previously) is then used by a system-level simulation framework to evaluate the implication of different memory technologies. These memory technologies are used for microprocessor caches at different levels. In [23], a cycle accurate performance simulator like gem5 is employed as using this simulator, it is possible to configure suitable cache parameters such as capacity, associativity, latency, block size and policy. Figure 2.24 shows the gem5 based system evaluation platform which is tuned to support both single core and multi-core capabilities. Apart from the cache parameters from NVSim, gem5



Figure 2.24: A system-level simulation framework to evaluate the implication of memories

(a) performance

(b) energy

(c) performance

(d) energy

Figure 2.25: System-level comparison of SRAM, STT-MRAM and SOT-MRAM with their nominal conditions [23].

simulator also requires system-level applications as input to execute the simulations. The gem5 tool generates the output for overall system performance (e.g. runtime) and cache statics such as read and write access per cache.

Using aforementioned framework, we have compared the SRAM, STT-MRAM and SOT-MRAM technologies in their nominal conditions. For that, we have used out-of-order processor with clock frequency of 3GHz and L1 (both data and instruction), L2 and L3 (shared) caches are configured as 32KBytes, 512KBytes and 16MBytes, respectively. They executed several applications such as MiBench, SPEC2000 and SPEC2006, in additional to several multi-core workloads and generated the final results as depicted in Figure 2.25.

## 2.6 Summary

In this chapter, first the classification of the semiconductor memories was discussed, in which the conventional memories such as SRAM, DRAM and flash, were explained in detail with their pros and cons. Then, the evaluation of spintronic technology was discussed in chronicle order. In that, initially, the discovery of *Giant Magnetoresistance* was discussed, latter *Magetic tunnel Junction* device and its associated *Tunnelling Magnetoresistance* was explained. The TMR is

essentially for the memory read operation, however, for the memory write operation, a controlling mechanism is required that can facilitate the magnetic switching for the free layer of the MTJ cell. For that purpose, several free layer switching mechanisms discussed in this chapter including, *Field Induced*, *Thermal Assisted*, *Spin Transfer Torque* and *Spin Orbit Torque*, and their associated shortcomings/advantages were also mentioned. Later in this chapter, MRAM characteristics discussed, in which general MRAM cell related behavior such as asymmetric switching, stochastic switching, read disturb, degradation effect, process variation, defects in MTJ's and fault modeling were explained. Afterwards, a flip-flop architecture implementation was explained, followed by a cross layer memory design approach starting from the bit-cell, all the way up to the complete memory architecture was explained with their system-level evaluations.

# 3 Self-timed read and write operations in STT-MRAM

To make STT-MRAM memory energy efficient, a self-timed technique for read and write operations [22] is proposed in this chapter. Using this technique, the energy is consumed only for the necessary durations for their respective memory operations. This is achieved by generating bitwise completion signals for read and write operations using a self-timed circuit that detects the end of a memory operation and initiates the closure of the corresponding read/write circuitry. By that means, this technique does not only improve the energy consumption but also the reliability. Since our technique tracks the actual end of a memory operation (read/write), it does not increase the read or *Write Error Rate* (WER), but it can also detect write errors and hence reduces the overhead for write error detection/correction mechanisms. In this chapter, the proposed self-time write technique is explained first, followed by the self-time read technique is described. Afterwards, a comprehensive result evaluation is done and finally, in the end of the chapter has the summary and conclusion section.

## 3.1 Overview

In STT-MRAM, when the Word-line (WL) is active, the current continuously flows through the MTJ cell and all components, including the peripheral circuitry, remain operational. Considering the fact that the active period of the WL is determined by the worst timing scenario, a very conservative timing guardband is chosen for read and write operations. This in turn results in a significant unnecessary current flowing through the bit-cells and hence a high dynamic power consumption. Moreover, a longer duration of the operational period also leads to reliability issues. For instance, it exacerbates the TDDB effect of the device during write operations and increases the read disturb rate for read operations [131].

In general, most of the write operations finish early enough, not (fully) exploiting the given timing guardband. For instance, in our design, the average switching time for 'AP' and 'P' is 9.7 and 3.9 ns, respectively. The write energy and WER for various write periods for both 'AP' switching and 'P' switching are shown in Background Section. If the write termination time could be accurately detected on the fly, the write circuitry can be closed earlier, resulting in less energy consumption, without affecting WER. In addition, due to the stochastic nature, there are also some write operations that may fail to complete within the given time period. If we would be able to detect the completion time, we can also detect these unsuccessful write operations and initiate an efficient error correction/retry mechanism.

To achieve this goal, some approaches were already proposed to detect or predict the completion of write operation in STT-MRAM. However, none of these solutions is comprehensive and can be used for all operating conditions. The early write termination technique targets only redundant write operations [132, 133], whereas verify-one-while-writing [134] and asynchronous asymmetrical write termination [135] only address the slow writes and fast writes, respectively. Furthermore, a variable-energy writes scheme [136] as well as a self-terminated write driver [24] approach were proposed that can tackle all write operations. However, these techniques employ an inverter (or buffer) as write completion detector which is impractical due to two reasons. First, it is hard to control the trip point of the inverter for the magnetic flip

Figure 3.1: Conceptual diagram to detect the switching of the MTJ cell

detection in the presence of variations, in particular, for smaller technology nodes. Second, there is a high short-circuit current through the comparator which significantly increases the total energy consumed during a write operation. Therefore, a robust bit flip detection technique is needed instead of the simple inverter- based comparator design. In contrast to all these write operation schemes, no technique is yet available that addresses read operations in STT-MRAM. However, various asynchronous circuit design approaches are implemented for other memory technologies [137–140]. Unfortunately, these circuit techniques are ineffective to track the MTJ-based bit-cell behavior. Hence, for STT-MRAM, an asynchronous technique which can track all memory operations in all operating conditions and terminate them on the fly is necessary for efficient energy reduction.

For a reliable and low-power implementation of STT-MRAM, it is mandatory to have a technique that can detect the actual end of all types of memory operations and closedown the corresponding circuitry directly after the operation ends. Therefore, we propose a self-timed bitwise approach that allows an on the fly, asynchronous termination of read and write operations as soon as these are completed. This technique not only reduces the power consumption, as motivated before, but also improves the reliability as current is flowing only in considerably reduced time periods.

## 3.2 Self-timed Technique

In this section, we explain our proposed methodology in detail. Section 3.2.1 and Section 3.2.2 explain the proposed circuit-level technique for an asynchronous detection and termination of write and read operations, respectively.

### 3.2.1 Asynchronous Write Termination

As explained previously, the write behavior in STT-MRAM is asymmetrical, i.e. writing '1' ('P') is much faster than writing '0' ('AP'). In addition, the write process is stochastic, i.e. the actual time required for a write operation is highly variable. This also applies to the bit-wise write operations inside the same data word. Hence, to achieve an acceptable WER usually very conservative timing windows are employed. However, improvement in WER has a trade-off with a high energy consumption. To address this problem, here we propose a self-timed based bit-wise write termination technique. This technique can detect the actual end of a write operation using a detection circuit based on a current sensor, and closes the corresponding write

circuits afterwards to reduce the power consumption. This behavior is explained below in two phases, *write detection* and *write termination.*

### Phase I: Write Detection

Since the bit-cell state in STT-MRAM corresponds to a resistance value, a write operation that changes the bit-cell value ('0'→'1', or '1'→'0') also changes the bit-cell resistance. Consequently, the write current changes accordingly. Based on our implementation, this current difference is around $90\,\mu$A (see Figure 3.1), which is enough to be used by our technique to detect the end of a write operation.

In fact, while this current difference itself can be easily detected, the final detection circuitry has to deal with two main challenges: 1) The write current has a bidirectional path, i.e. the current direction for writing '1' and writing '0' are opposing and different in magnitude. 2) Non-differential write operations ('0'→'0', or '1'→'1') do not lead to a detectable current change. To deal with these challenges, we developed a detection circuit as shown in Figure 3.2. As one can see, there are two different branches in the circuit, one is used for writing '1' (*P-Branch*) and the other one for writing '0' (*AP-Branch*). Having these two different branches, the bidirectional current path is no longer an issue, since there is a separate branch for each write operation (and hence current direction). Furthermore, this implementation also allows to detect non-differential write operations, as we will explain later. The output of these branches (*out_AP* and *out_P*) can then be used to terminate the write operation. To reduce the power overhead



Figure 3.2: Circuit diagram for asynchronous bit-wise write termination (AP-branch and P-branch are added for our technique)

Figure 3.3: Circuit diagram for the Activation Control Circuits (ACCs) of the 'P'- and 'AP'-branch of the bit-wise write detection technique

of this scheme, depending on the write operation, only the required branch is activated (via write_{AP,P}_enable, WL, out_{P,AP}).

Each branch contains its own sensing circuit. These circuits are intended to detect a current difference between a reference cell and the bit-cell. As soon as this current difference is around a particular threshold (90 $\mu$A in our setup), the ongoing write operation is finished and the sensing circuitry can trigger a signal (out_{P,AP}) to initiate the write termination. Therefore, each sensing circuit requires two inputs. The first input is biased through the write current (bit-cell) using a current mirror and the other input is biased through a reference cell. Since the current direction and difference need to be known in each branch, the reference cell has to be designed accordingly. Hence, if a write '0' ('AP') is going to be tracked, the reference cell has to have a 'P' configuration to achieve a current difference between the bit-cell and the reference cell. For writing '1', it is the other way around. These sensing circuits are enabled using sen_en_{AP,P} signals which are controlled using *Activation Control Circuits* (ACC). The ACC consists of a 3 input AND gate and generates the enable signals (sen_en_{AP,P}) to turn ON/OFF the required sensing circuits using tail transistors N13/N14. ACC circuits have three inputs namely, WL, input data and out_{AP,P}. The out_{AP,P} signal is fed back to the ACC signal to deactivate the sensing circuits once when switching is over. Consequently, at the end of all kinds of write operations (differential and non-differential) a current difference between the reference cell and bit-cell will exist. As this current difference (of around 90 $\mu$A in our setup) is sensed, an acknowledgement signal (out_{P,AP}) is triggered to indicate the end of the actual write operation. Note that in the case of non-differential write operations the bit-cell current remains constant and is in its final state right from the beginning. Hence, the acknowledgement signal is generated almost immediately. In our implementation, the delay of the acknowledgement signal with respect to WL is around 198 ps.

A limitation of the implementation explained above is that in the case of differential write operations the activated branch is active for a very long period of time. This is because the selected branch is enabled with the activation of WL and remains active until the flip is observed. Consequently, the static current flowing through the CMOS gates leads to a high power consumption. To make our approach power efficient, it is required to enable the selected branch (sensing circuit) only when it is needed. In other words, the sensing circuit should be active only in a short time period during which the flip occurs. Therefore, the ACCs are modified as depicted in Figure 3.3. As illustrated in Figure 3.3, each ACC consists of two branches, one for differential write operations and one for non-differential operations. Both branches contain a delay element to activate/deactivate the detection circuitry at the right point in time. While the non-differential branch relies on a simple inverter chain to achieve the required active period ($\approx$ 198 ps), the differential part uses a clock based delay approach as explained in [135] to predict the point in time at which the flip occurs. As a result of these two

Figure 3.4: Sequence of ACCs outputs and write detection signals (1 = First activation period of the sensing circuitry to detect non-differential writes, 2 = Second activation period if write is $AP \to P$, 3 = Second activation period if write is $P \to AP$)

branches, there is always an active period for the detection circuitry right at the beginning of the write operation to detect non-differential operations as shown in Figure 3.4. This period is inevitable, even if differential write operations occur, as it is impossible to distinguish between differential and non-differential write operations before the operations start. Hence, in the case of differential write operations, there are two active periods (one right at the beginning, and another one for the differential write operation). The start time of the second active period depends on the actual write operation, i.e. for $AP \to P$ the second active period is around 3.9 ns, while it is around 9.7 ns for $P \to AP$ operations (see Figure 3.4).

**Optimal Write Termination Point**: Due to the stochastic nature, if write operations finish much earlier, this will not be detected until the corresponding detection branch is activated. Hence, to reduce the power overhead of the detection circuitry, a certain amount of



Figure 3.5: Optimal write termination point (for a 32 bits word) to enable the write detection circuit

energy savings due to a "faster than expected write" is sacrificed. In order to obtain an optimal write termination time to enable the write detection circuit using clock based delay signal, we extracted the energy consumption using SPICE simulation for the write distributions as described earlier in Background Section. The energy consumption for various termination points for 32 bits is as shown in Figure 3.5. The figure depicts that the optimal termination point for 'AP' switching is around 10 ns where the energy consumption is minimum. Similarly, the optimal termination point for 'P' switching is obtained as 4 ns.

### Phase II: Write Termination

Once the end of the write operation is detected, the acknowledgement signal (out_{AP,P}) makes a transition from '1' to '0'. This can be used to terminate the write operation asynchronously. Therefore, the acknowledgement signal is fed back to the ACCs as input to deactivate the branches, and it is also employed to cut off the power for the bit-cell, similar to the method mentioned in [135]. Accordingly, the current flow through the bit-cell is stopped. As this is a *bit-wise* termination, for each bit-cell the write operation is terminated at the optimal point in time, which minimizes the overall current flowing through the bit-cells. This in turn, will improve the overall memory reliability, due to reduced TDDB, as the amount of current passing through the MTJ and the oxide barrier is minimized.

Note that the clock based delay signal [135] of the ACCs cannot be used directly to terminate the write operation. This is because the write latency of the actual bit-cell(s) can be significantly different from the predicated delay due to variations and in particular, the stochastic write behavior. Hence, this technique can be used only for power reduction. Furthermore, as soon as all bit-wise write operations are terminated (i.e. for all bits the out_{P,AP} is '0') a combined acknowledgement signal is generated to deactivate the WL signal. It is noteworthy that the write latency is still determined (asynchronously) by the slowest bit-cell, however, write circuits are disabled bit-wise, hence only necessary energy per bit-cell is used.

**Write Error Detection**: A major problem with write operations in STT-MRAM is their stochastic nature, i.e. the actual time which is required to perform the write operation is non-deterministic. Hence, even with conservative timing margins, it may happen that the write operation is not yet finished, i.e. when the word-line is deactivated. In these situations the acknowledgement signal (out_{P,AP}) is still '1', but the word-line signal is already '0'. Hence, with this knowledge an error signal (WE_detect) can be created according to:

$$\text{WE\_detect} = \text{out\_\{P,AP\}} \wedge \overline{\text{WL}}.$$

This signal can be used by the memory controller to trigger an error correction mechanism (e.g. retry the write operation). By this means, write errors can be detected very early (immediately after they occurred), and very fast as well as efficient error correction mechanisms can be invoked. As a result, this allows to reduce the overhead of sophisticated *Error Detection Code* or ECC techniques.

### 3.2.2 Asynchronous Read Termination

To generate the read completion signal, we propose a circuit-level technique added onto the existing read methodology. In our proposed technique, the logic values, which are complementary to each other, at the output ports of the sensing circuit are exploited dynamically when a read operation is performed. Therefore, for each bit, the actual end of the read operation is detected on-the-fly and the read completion signal is generated accordingly. When the read completion signals for all bits are generated, meaning the read operation for an entire word is accomplished, a combined acknowledgement signal is generated to deactivate the read process

Figure 3.6: Circuit diagram for a bit-wise read detection technique (dotted boxes are existing circuits, shaded parts are added for our technique)

of the memory. Please note that, unlike the asynchronous write termination technique, in this technique we perform word-wise termination instead of the bit-wise as both read '0' and '1' have similar delay values. Hence, the entire process of the asynchronous read termination can be split into two phases, *i.e. read detection* and *read termination*, that will be explained next.

### Phase I: Read Detection

For the implementation of the read circuitry in our technique, we have chosen a fast and low-power pre-charge sense amplifier similar to the one presented in [141]. This sense amplifier, as shown in Figure 3.6, is equipped with an *equalizer* to keep the potential of the two output nodes ($q1$ and $q2$) equal. The sense amplifier consists of two branches, namely, *bit-cell* and *reference cell branch*. In the bit-cell branch, all the bit-cells of a single column are connected, however, out of them, only one is activated using word-line (WL) signal based on the address bits of the memory. The reference branch consists of four MTJ cells connected in such a way that the effective read resistance is the middle of the two possible resistance states ($R_{AP}$ and $R_P$), i.e. $(R_P + R_{AP})/2$ where $R_P$ and $R_{AP}$ are the resistances during 'P' and 'AP' magnetization states, respectively. Similar to the bit-cell branch, several bit-cells (equal in number to those of the bit-cell branch in a column) are connected to this reference branch in order to maintain the same load effect on both branches. The access transistor of the reference cell is driven using an *Address Transition Detection* (ATD) circuitry [142].

The chosen sense amplifier has two phases of operation:

- Pre-charge phase: The *pre-charge* (PC) signal is '0' and the equalizer circuit is ON which keeps $q1 = q2 = 1$.

- Evaluation phase: The PC signal is '1' and the equalizer circuit is disabled. Due to this, the potential of the two nodes $q1$ and $q2$ starts gradually to decrease. However, the

Figure 3.7: Simulation waveform for a bit-wise read detection

potential of the node which has less resistance decreases faster than that of the other node, which increases the potential of the other node due to back-to-back connected inverters (see Figure 3.7). As a result, either of the two output nodes ($q1$ and $q2$) stables at '0' and the other remains '1' (depending on the resistance state of the two branches). Whenever these potentials are stabilized, the value is read from the bit-cell.

In the pre-charge phase, both $q1$ and $q2$ nodes have the logic value '1', while after the completion of the read operation (i.e. bit-cell value successfully read) either of them is '0'. This behavior can be detected with an AND gate as shown in Figure 3.6. Hence, the output of this gate is '1' at the beginning of a read operation (during the pre-charge phase) and '0' as soon as the data is read (end of the operation). Accordingly, this output signal indicates the end of the read operation for a particular bit-cell. Since, these output nodes have low impedance, the AND gate can be easily driven. This behavior is illustrated in the waveform as shown in Figure 3.7. Please note that, the required AND gate is not part of the critical path, and hence does not affect the overall read latency. Furthermore, our technique also works with any other sense amplifier such as [141, 143, 144].

### Phase II: Read Termination

As soon as the read acknowledgement signals for all bits in a row are generated, a combined signal is generated to deactivate the word-line. This is achieved using a NOR gate as shown in Figure 3.6. Once the word-line is deactivated, the bit-cell is isolated from the sense amplifier and the memory is ready for the next operation.

## 3.3 Experimental Setup and Results

In order to evaluate the effectiveness of our proposed technique for both read and write operations, we implemented the proposed technique at the circuit-level in SPICE environment. The circuit-level analysis is then projected at the architecture-level to evaluate the effect of the proposed technique on the energy consumption of various workloads. Please note that, we have used an in-plane STT model for our implementation because this one is validated with fabrication process, however, our proposed self-timed technique can also work efficiently with perpendicular STT model.

Table 3.1: MTJ parameters

| Parameters | Value |
|---|---|
| 'AP' resistance | 6.5 KΩ |
| 'P' resistance | 2.6 KΩ |
| Damping factor | 0.2 |
| Size of MTJ | 40 nm × 90 nm |
| Read current | 12 μA |
| Write critical current | 51 μA |

### 3.3.1 Circuit-level Analysis

For the circuit-level implementation, we employed TSMC 65 nm general purpose SPICE models and the MTJ model presented in [108]. The MTJ parameters used in this work are summarized in Table 7.1. The simulations are conducted with the supply voltage of 1.2 V and temperature of 27°C.

Since the proposed technique reduces the energy consumed after the completion of the desired operation until the end of the clock period, its energy saving and reliability improvement are a function of the clock period. This fact will be evaluated next, with the baseline of standard STT-MRAM design.

**Evaluation of the Read Methodology**

Figure 3.8 shows the percentage of the reduction in energy and read disturb rate for the read operation across different clock periods ranging from 1 ns to 23 ns. As one can see, the energy saving of read operation grows by increasing the clock period, which is due to the fact that the sense amplifier requires less than 300 ps to read a single bit-cell. Consequently, the lowest energy savings (around 30 %) are obtained for aggressive clock periods, while for low clock periods the read energy savings can even reach 92 %. Similarly, the reduction in the read disturb rate shows a saturation behavior. As the read disturb probability is time duration dependent, the read disturb rate for our proposed technique can be reduced by more than 60 % for a clock frequency of 1 GHz. Large clock periods improve this result up to 98 %.



Figure 3.8: Read energy saving and read disturb reduction with respect to different clock periods

Figure 3.9: Write energy saving for all possible switching cases with respect to different clock periods

## Evaluation of the Write Methodology

Similar to the evaluation of the efficiency of our proposed read methodology, Figure 3.9 shows the analysis of our proposed write approach for all possible write operations (differential and non-differential) reflecting the mean write delay. Since a write access requires more than 1 ns, it might require multiple clock cycles to accomplish. This results in a "zig-zag-behavior" as shown in Figure 3.9. Whenever, the write delay is matching with a multiple of the clock period, the word-line is deactivated right after the write operations finish. Hence, in these cases the savings of our proposed technique are worst, since it does not allow to terminate the write operation earlier, but adds some energy overhead due to the required circuitry. The power consumption due to the write detection circuit for 'P' and 'AP' switching is 565 $\mu$W and 304 $\mu$W, respectively. If all write operations have similar occurrence rates (i.e. 25 %), our proposed write scheme provides more than 70 % energy saving in average. In contrast, if the write delay is much more than a multiple of the clock period, the energy savings are much higher, and grow with the increase in the clock period.

## Area Analysis

We have estimated the area overhead of our proposed self-timed read and write technique. The total memory area is obtained for a 2 MByte capacity using the NVSim tool [129]. The area of additional gates, which are required for every write circuitry of our proposed technique, are obtained from the TSMC standard cell library with appropriate drive strengths. The area of the write detection sensor circuits is estimated based on these transistor sizes. In addition, we included the area of a 4-bit counter and the write termination circuitry. With all this, the total area overhead of our proposed technique is obtained as less than 0.4%.

## Process Corner Analysis

In addition to the design-level analysis, we also performed a process variation analysis at circuit-level. For that, we first extracted the delay, current and power values for a single bit-cell using SPICE simulations for three process corners, i.e. *fast*, *typical* and *slow*. Please note that we have considered process variation effects separately for the MTJ and CMOS devices, as these two are different fabrication technologies. For the MTJ cell, we have considered $\pm 3\sigma$ variations for the *Tunneling Magneto-Resistance* and the product of *Resistance and Area* values. After extracting

Figure 3.10: Write energy saving for write '1' and write '0' for various memory sizes for our proposed write technique

the cell-level information, we fed the corresponding corner case values separately to the NVSim tool [129] in order to obtain the energy savings at memory architecture level. The write energy savings for our proposed technique for various memory sizes with the baseline of standard STT-MRAM design, are shown in Figure 3.10. In this figure, the write '1' ('AP'→'P') energy savings are more compared to that of the write '0' ('P'→'AP'). This is because, the average switching latency for write '1' is much shorter than the write '0' and hence, the write termination for write '1' happens much earlier, which results in higher energy savings. Moreover, write '0' energy savings have more variations than that of the write '1' energy.

In a similar way, we explored the process variation effects for the self-timed read technique as well, and we found that these variations are not that pronounced in our case. The total read time for a single bit-cell has the best and worst case latencies of 333 ps and 345 ps, respectively. Because of this fact, we employed a word-wise approach for the read termination as mentioned in Section 3.2.2. Nevertheless, if process variation is more severe, a bit-wise termination can also be chosen.

### 3.3.2 Comparison with state-of-the-art

To demonstrate the merits of our proposed technique, we have compared it with the existing self-terminated write technique [24]. To do so, we implemented this technique in our simulation framework. In this technique, an inverter based comparator is used as detection circuit which is the key element of the design. For such a comparator design, one of the transistors (either NMOS or PMOS) operates in the active region and the other one in the saturation region resulting in high short-circuit power consumption. However, the detection circuit for our proposed technique is implemented based on a reliable current sensing technique which is enabled only for the required duration when it is necessary (see Section 3.2.1). Due to this optimization for the active duration of the detection circuit, it is overall energy efficient. The energy and detection delay values for the existing technique along with our proposed technique are given in Table 3.2. In order to account for the stochastic effect, these values are obtained with respect to a write period of 18 ns.

Beside comparison with existing technique, we also performed a memory architecture-level comparison with SRAM using the NVSim tool [129]. This comparison for a 2 MByte memory capacity for a word size of 512 Bits is summarized in Table 3.3. As shown in the table, STT-MRAM has significant advantages over SRAM in terms of area, read latency and leakage

Table 3.2: Comparison of our proposed self-timed technique with a self-terminated write driver technique [24] for a single bit-cell

| Parameters | | Self-terminated write driver [24] | Self-timed technique |
|---|---|---|---|
| Write Energy [pJ] | 'P'→'AP' | 3.8 | 1.2 |
| | 'AP'→'P' | 2.8 | 0.3 |
| Write detection delay [ps] | | 246.0 | 198.0 |
| Write detection using | | Inverter based comparator | Current sensing |
| Optimization to enable the write detection circuit | | No | Yes |
| Type of memory operation | | Write | Write and read |

Table 3.3: Comparison of SRAM, standard STT-MRAM and proposed self-timed read and write technique for a 2 MB memory

| Parameters | SRAM | Standard STT-MRAM | Self-timed Read-Write Technique |
|---|---|---|---|
| Area [mm2] | 11.24 | 5.96 | 5.97 |
| Read latency [ns] | 1.35 | 1.35 | 1.11 (nominal) |
| Write latency [ns] | 1.16 | 18.73 | 10.74 (nominal) |
| Read energy [nJ] | 1.19 | 0.61 | 0.49 |
| Write energy [nJ] | 0.76 | 1.26 | 0.51 |
| Leakage [W] | 3.90 | 1.55 | 1.55 |

power. In contrast, SRAM is better for write latency and both read and write energies. As mentioned earlier and also evident from the Table 3.3, our proposed self-timed technique has an improvement in terms of both read and write energies over the standard STT-MRAM design. Moreover, in standard STT-MRAM design, the latency values are always set according to the worst-case, however, using our proposed technique the average latency improves as shown in the table.

However, it has some overheads for the area ($< 0.4\,\%$) and leakage power ($< 0.1\,\%$) due to the additional circuitries added to the standard design.

### 3.3.3 Architecture-level Analysis

In order to show the architecture-level implications of the proposed approach, we conducted an experiment to evaluate the energy saving with respect to realistic memory accesses for different workloads. In this regard, we investigated the energy saving for a 32 MByte main memory of Leon3 System-on-Chip (SoC) as a case study. This SoC has a minimum clock period of 1.24 ns, hence, the main memory is implemented with the STT-MRAM model considering the read and write operations of 6 and 19 cycles, respectively.

The HDL description of the main memory is modified to compute the read and write statistics as well as the probability of having different write types ('AP'→'AP', 'AP'→'P', ...). Several workloads from MiBench suite [145] are executed on the Leon3 processor and the energy savings are computed based on the operation types for their respective workloads. Table 3.4 summarizes the results for the architecture-level analysis. It is shown in the table that the proposed approach reduces the overall energy of read and write operation for the main memory by 66 % and 89 %, respectively. Also, the overall energy saving is 88% which is almost identical to the write energy saving as the overall absolute energy consumption of the write operation is very high compared to that of the read operation.

Table 3.4: Energy reduction for leon3 processor for main memory for various MiBench workloads

| Benchmark | Cycles | Write Type Occurrence [%] | | | | Read Occurrence | Energy Saving | | |
|---|---|---|---|---|---|---|---|---|---|
| | | AP→AP | AP→P | P→AP | P→P | | Read | Write | Overall |
| stringsearch_small | 500M+ | 43.17% | 0.43% | 0.26% | 1.21% | 54.93% | 66.11% | 91.77% | 90.96% |
| stringsearch_large | 77M | 40.07% | 0.95% | 0.95% | 1.35% | 56.68% | 66.17% | 90.40% | 89.59% |
| basicmath_small | 500M+ | 18.96% | 3.42% | 3.49% | 6.05% | 68.09% | 66.08% | 82.95% | 82.10% |
| bitcount_small | 500M+ | 45.62% | 0.26% | 0.29% | 0.49% | 53.35% | 66.16% | 91.82% | 91.05% |
| crc32 | 6M | 29.74% | 0.11% | 0.11% | 0.29% | 69.75% | 65.97% | 92.03% | 90.52% |
| qsort_large | 3M | 30.85% | 6.19% | 3.83% | 1.01% | 58.11% | 66.17% | 84.33% | 83.71% |
| Average | | 33.73% | 1.89% | 1.49% | 1.73% | 60.15% | 66.11% | 88.88% | **87.99**% |

### 3.3.4 Time Dependent Dielectric Breakdown Analysis

In addition to the design analysis in terms of energy and area for our proposed techniques, we also performed a reliability analysis considering *Time Dependent Dielectric Breakdown* (TDDB). In our proposed self-timed based write scheme, the current value remains same, however, the duration of the current flow is reduced significantly.

The bit-cell current value is obtained using SPICE simulations, and the duration of the current flowing is considered based on the average switching delay of the bit-cell. Using the equation for TDDB (as described in Background) and the architecture-level results for the write type occurrence for various workloads, the average TDDB improvement is obtained as over 90 %. This is because, for a large number of non-differential writes ('AP'→'AP' and 'P'→'P'), our proposed technique terminates the write operation right at the beginning, which means almost no current flowing through the bit-cell. Moreover, the fast 'AP'→'P' transitions occur more often than the slower 'P'→'AP' operations. Overall, with our proposed self-timed technique, the write current flows through the bit-cell only for the required duration, hence such a high TDDB improvement.

## 3.4 Conclusions

STT-MRAM is an emerging non-volatile memory technology due to its various beneficial features such as scalability, high density, no leakage and high endurance. However, high dynamic power is still a major concern for this memory technology. To overcome this, it is required that all the active components are turned-off immediately after they finish their respective operations. We proposed a self-timed bit-wise termination technique for both read and write operations in which the operation completion is detected on-the-fly and an acknowledgement signal is generated, with which the respective active components can be turned-off. Our results for Leon3 main memory show that the proposed technique can reduce the overall memory energy consumption by 88 %.

# 4 Improving Write Performance for STT-MRAM

In this chapter, we propose two complementary bit-wise write schemes for STT-MRAM, a static and a dynamic technique [146]. In the static technique, the latency of the slow-write is reduced by addressing the bit-cell asymmetry. Therefore, the bit-cell structure is modified and an increased current is passed through the bit-cell only for the slow-write operations. The dynamic technique reduces the write margin by further increasing the write current in a step-wise manner until the bit-cell switching process is completed. Moreover, the write margin due to process variation can also be reduced using our proposed dynamic technique. This chapter is organized as follows: related work is discussed first, then our proposed static and the dynamic write techniques are explained. Afterwards, comprehensive results are presented, followed by the summary and conclusion of this chapter is mentioned.

## 4.1 Related work

The write latency of STT-MRAM can be reduced at device-level by relaxing the thermal instability factor [147] at the cost of increased retention failures [148–150]. Although this technique reduces the write latency, it does not address the write asymmetry issue. In order to address this issue, the technique presented in [151] balances the write delays by changing the voltage-level of the memory word line based on the input data. By applying this technique, the effective latency becomes the average of the two write latencies. This technique not only significantly increases the fast-write latency, but also requires two level converters and an additional negative power supply which eventually increases the total energy. The write latency and energy of the STT-MRAM are significantly reduced in [152], however, an additional magnetic field is used to assist the switching process. Furthermore, an adaptive write current boosting technique is proposed in [25] that addresses the process variation induced slow bit-cells. Since the write current is boosted for the entire column, this technique consumes high energy for the fast switching bit-cells in that column.

In summary, a suitable solution for balancing the write latency at low cost is still missing and requires more attention. Furthermore, none of the existing techniques has explicitly addressed the write margin due to the stochastic write behavior, and thus none of them exploited this important knob to improve the write latency and energy. Therefore, it is important to effectively reduce both the write margin and the slow-write latency in order to shorten the overall write latency.

## 4.2 Proposed Write Latency Reduction Technique

### 4.2.1 Overview

As explained earlier, the write access latency in STT-MRAM depends on the slow-write latency considering appropriate amount of timing margins for the stochastic and process variation effects. There are several techniques trying to optimize the write current to improve either yield [153], reliability [134, 154], or energy efficiency [24, 133, 136]. However, there is no

Figure 4.1: Overview of the proposed static and dynamic write technique.



Figure 4.2: Current versus delay, for both 'P'→'AP' and 'AP'→'P' switching.

work targeting the latency of slow-writes and the associated write margins through the current optimization. Hence, we propose bit-wise static and dynamic write techniques to reduce the total write period, as described in the flowchart shown in Figure 4.1. In the static technique, we increase the write current for the slow-writes in such a way that the delay of the slow-writes ('AP' switching) becomes equal to that of the fast-writes ('P' switching). Then both 'AP' and 'P' transitions are monitored per bit-cell using a switch detection scheme [24]. If the desired switching is incomplete after a given duration of time, the write current is increased for that bit-cell by a certain amount for a specified duration. In case that bit-cell switching is still incomplete, the write current is further increased and these steps are repeated until the corresponding bit-cell switching is completed.

## 4.2.2 Current versus delay in MTJ

In STT-MRAM, the switching delay depends on the amount of current flowing through the MTJ cell, i.e. a larger current value leads to a shorter switching delay. The current and delay relations for both switching types are depicted in Figure 4.2. These values are obtained using the in-plane based MTJ device model presented in [108]. The switching delay is obtained for a single MTJ cell by sweeping the current value. As shown, the slow switching type ('P'→'AP') varies significantly for a slight increase in the current. On the contrary, the delay of the fast switching type ('AP'→'P'), decreases on a small scale even with a considerable increase in the

Figure 4.3: Circuit diagram for the implementation of the static technique (blue transistors and gates are added to the default design). Transistors MN3 (W=600 nm) and MP3 (W=1200 nm) are introduced to increase the 'AP' switching current.

current. Therefore, it is beneficial to target the slower operations for an overall reduction of the write period.

In order to maximize the gain for both performance and energy, the 'P'→'AP' switching latency should closely match with that of the 'AP'→'P' switching. Therefore, the write current for the slow-writes can be increased to the extent that the switching delay of the slow-writes matches with that of the fast-writes. For this configuration, almost the same latency (around 3.9 ns) can be achieved for the fast- and slow-writes using the current values of 138 uA and 180 uA, respectively.

### 4.2.3 Static write technique

In the static write technique, our intention is to shorten the slow-write latency by reducing the asymmetry of the bit-cell. The MTJ based asymmetry is handled by increasing the write current only for the slow-writes. For this purpose, our proposed write circuitry is designed in such a way that it drives an increased current for the slow-write operations. However, due to the voltage drop across the access transistor in the 1T1MTJ bit-cell structure, the current increase is relatively small. Therefore, to overcome this limitation, a PMOS transistor is introduced in the bit-cell.

The circuit diagram for the implementation of the proposed static technique is shown in Figure 4.3. As shown in the figure, a PMOS is connected in parallel to the original NMOS access transistor (i.e. a transmission gate structure) in order to eliminate the voltage drop due to the NMOS transistor. These PMOS and NMOS transistors are conditionally used for memory read and write operations. During the write operation, both transistors are ON, so that write currents in both directions flow evenly without any degradation, whereas only the PMOS is ON during the read operation. A low current is desirable for the read operation [119, 155, 156] and since PMOS has lower mobility, it delivers less current. Note that, with the proposed technique, the read and write currents can be optimized separately. The activation conditions for these two transistors are obtained using an AND gate which has two inputs, namely, *write enable* (WE) and *word line* (WL). WE is considered as HIGH and LOW for the write and read operations, respectively. The WL signal is driven through the address decoder

Figure 4.4: Waveform for the static method with balanced 'P'→'AP' and 'AP'→'P' switching.

of the STT-MRAM. The output of the AND gate is connected to the gate terminal of the NMOS and WL is directly connected to the gate terminal of the PMOS. Hence, the PMOS is always ON whenever WL is activated irrespective of the type of the memory operation. On the other hand, the NMOS is ON only during the write operation, i.e. when WE is HIGH.

Furthermore, as illustrated in Figure 4.3, another important component of the static write technique schematic is the write circuit. It consists of two inputs, one is *input data* which is clock dependent and driven through the input latch, and the other one is the (WE) signal which is applied through the memory ports. The output of the write circuit drives the bit-line (BL) and the source-line (SL) terminals of a single column. These output values are always complementary to each other during the write operations. In order to demonstrate the two write current paths, the last pair of inverters is shown as transistor-level schematics in Figure 4.3. These inverters drive the BL and the SL. To write a 'P' configuration, the input data is HIGH which makes $\overline{BL}$ and $\overline{SL}$ as HIGH and LOW, respectively. As a result, the write current flows from BL to SL because the transistors MP1 and MN2 of the write circuit are ON while MP2 and MN1 are OFF. On the other hand, to write an 'AP' configuration, the input data is LOW which makes the transistors MP2 and MN1 ON, whereas MP1 and MN2 are OFF. In this case the write current flows from SL to BL to write the 'AP' configuration. Therefore, MP1 and MN2 are the driving transistors to write the 'P' configuration whereas MP2 and MN1 are the driving transistors to write the 'AP' configuration (slow-writes). In order to fulfill the given objective, i.e. to pass an increased current during the slow-write operations, it is required to strengthen the MP2 and MN1 transistors. This can be done either by up-sizing the MN1 and MP2 transistors or by adding a PMOS and an NMOS transistors as shown in Figure 4.3 In the actual implementation, these additional transistors can be included as layout finger and share the same diffusion.

The waveforms for both write operations are depicted in Figure 4.4. It shows that the magnetic switching latency of the two writes are closely matching (around 3.9 ns). Hence, using this static technique, both the inherent torque asymmetry of the MTJ cell and access transistor induced asymmetry are eliminated.

Figure 4.5: Concept of dynamic write technique.

## 4.2.4 Dynamic write technique

The stochastic switching behavior in STT-MRAM means that the actual write period to a bit-cell is random. Due to this behavior, the write operations for some bit-cells in a row of the memory array (i.e. bits in a word) are already completed while for others the write operations remain incomplete in the specified time period. These unfinished write operations contribute to WER. To account for those bit-cells whose transitions are not completed for a given time duration, an additional margin, that is almost as large as the average switching delay, is required to achieve the target WER ($10^{-9}$). By increasing the amount of write current, the probability of unfinished operations drops, and hence a smaller margin is required (see Equation 2.13). However, the write current cannot be increased too much due to energy and reliability constraints. Therefore, we propose a dynamic write technique in which the write current is increased in a step-wise manner only for those bit-cells whose write transition is not completed after a specified duration. This technique is implemented for each write circuitry (means bit-wise operation) in order to make it energy efficient.

In our proposed methodology, some parallel transistors, which are part of the write circuitry, are activated in a step-wise manner to drive more current, while the write operation is ongoing. The implementation is illustrated using a conceptual diagram as shown in Figure 4.5. The shaded region in this figure represents the average switching delay, therefore it is termed as *nominal region*. For the bit-cells with incomplete write operations (transitions) in the nominal region, the current through the bit-cell is dynamically increased. This is done at every edge of an available high frequency processor clock, which is a fraction of the write period of the STT-MRAM. For this purpose, an acknowledgement signal is required that can detect the completion of a write operation. This acknowledgement signal is generated through a circuit implemented in [22] as a switching detection signal. So, the write current that pass through the bit-cells with incomplete transitions, is increased in a step-wise manner (for every edge of available processor clock) until an acknowledgement signal is activated, i.e. the transition is completed. With this technique, the required write margin for the same target WER can be reduced (see Figure 4.9 for details). Please note that, using this technique, either the write period can be reduced for a given WER or the WER can be improved for the same write period. In this work, we choose the former one to improve the overall latency of the STT-MRAM.

The circuit-level implementation of the proposed dynamic write technique is shown in Figure 4.6(a). Here the write circuitry is the same as the one used for the static write technique. In order to increase the write current, several transistors are employed that are activated using conditional circuits ('C1' to 'Cn') based on the following conditions:

- activated on each edge of the external clock after the nominal region.

(a) Circuit diagram for the dynamic write. Several PMOS (W = 1200 nm) and NMOS (W = 600 nm) transistors are activated on conditional basis to boost the write current.



(b) Circuit diagram of the 'C1' component.

Figure 4.6: Circuit diagram for the implementation of the dynamic write technique.

- operated only when acknowledgement signal is not activated.

- activated for the respective write values, which is based on the direction of the current.

One of the conditional circuits, namely 'C1', is shown in Figure 4.6(b). The other conditional circuits are implemented similarly. This circuit is activated based on the aforementioned conditions and generates four outputs, from 'C1_out1' to 'C1_out4'. Here the direction of the current which is obtained from $\overline{BL}$ and $\overline{SL}$, decides whether it is an 'AP' or 'P' type of write operation and the corresponding output is activated for a bit-cell. Note that the current increasing steps can be independently optimized for 'P'→'AP' and 'AP'→'P' types of switching.

Table 4.1: Truth table for 'C1' circuit

| $\overline{BL}$ | $\overline{SL}$ | $\overline{Q}$ | ack_signal | c1_out1 | c1_out2 | c1_out3 | c1_out4 |
|---|---|---|---|---|---|---|---|
| X | X | X | 1 | 1 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 |
| 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |

We considered the acknowledgement as an active HIGH signal, meaning that this circuit will not activate any of the outputs if the acknowledgement signal is HIGH. Furthermore, to make it clock dependent, a D-Latch is used. The functionality of the implementation for the circuit 'C1' is explained using Table 4.1. The positive-level latches are used for the conditional circuit 'C1' and 'C3', while negative-level latches for 'C2' and 'C4'. Please note that, initially these latches are reset and then activated when the majority of write operations has been completed (after the nominal region).

**Evaluation of the number of steps:** Finding an optimal number of necessary boosting steps is challenging. The reason is that the distribution of the switching latency alters with the change in the current value, and it is not trivial to find the total write period for the target WER value in such a case. Therefore, to deal with this scenario, we obtain the switching delay based on the accumulated charge transferred. This is because the charge is ultimately responsible for the bit-cell switching. In general, the charge is computed as:

$$Q = I \times t \tag{4.1}$$

where I is the current and t is the duration of the current flow. The accumulation effect of the step-wise increase in the current value is obtained by adding the charge values for each boosting stage. The terminating point is the point where the accumulated charge is equivalent to the charge required for a WER of $10^{-9}$. In this regard, the only thing that need to be determined are the number of steps required, the duration of the steps and the increase of the current value in each step. As mentioned previously, the duration of steps are clock frequency dependent, and thus the only available knobs are the number of boosting steps and the current value in each step. Note that the durations of all steps are the same except for the last one, because it depends on the write period termination value. For 'N' steps, the total charge accumulation



(a) N=1      (b) N=2

(c) N=3      (d) N=4

Figure 4.7: Step-wise current increase for different suitable delays

Figure 4.8: Energy consumption for various delays of different equidistant step sizes

can be formulated using Equation (4.1) as:

$$Q_{total} = Q_{base} + \sum_{i=1}^{N} t_{s_i} \left( \frac{I_{max} - I_{base}}{i} + I_{base} \right) \tag{4.2}$$

where $I_{base}$ is the initial current value, $I_{max}$ is the maximum current value, $Q_{base}$ is the charge transferred during initial period (for nominal switching) and $t_{s_i}$ is the duration of the $i^{th}$ step.

Following this optimization methodology we obtained the boosting patterns depicted in Figure 4.7 for a clock period of 500 ps and a maximum current of 230 $\mu A$. As it can be seen, for all evaluated values of 'N', the total write period is almost the same. This is due to the fact that the delay sensitivity on the current is rather small, when the write current is larger than 180 $\mu A$. Therefore, we choose the solution with the lowest energy consumption. The energy consumption is computed for each step of current increase for the given duration of the time, as illustrated in Figure 4.8. This figure shows that the lowest energy consumption is achieved for 'N=3' (with $t_s = 0.5$ ns) and 'N=2' (with $t_s = 1.0$ ns). Thus, we select 'N=2' for our implementation as it requires a simpler boosting circuitry compared to 'N=3'.

### 4.2.5 Discussion

In the static write technique, the two switching delays are balanced with the 'AP' current value of 180 uA, for a particular experimental setup. For that we have increased the width of the existing access transistor (from 470 nm to 600 nm) in addition to the introduction of the PMOS transistor (size of 600 nm). For the dynamic write technique, the size of both access transistors are adjusted (as 770 nm) to be able to deliver the final current value that is 230 $\mu A$. Please note that the sizing of the transistors will change, if perpendicular STT is used instead of in-plane STT, but the overall methodology remains the same. Furthermore, such an increase in access transistor sizes and the introduction of several other devices result in an increase in the leakage power in the write circuit (by 3.36 nW per write circuit). This leakage overhead is negligible compared to the total leakage power of the memory (193 mW for 512 KByte capacity). Please note that in STT-MRAM, only periphery circuits contribute to the leakage power.

## 4.3 Experimental Setup and Results

To evaluate the effectiveness of our proposed techniques, we analyzed the circuit-level implementation first and then projected the obtained results to an architecture-level evaluation.

Table 4.2: MTJ parameters

| Parameters | Value |
|---|---|
| 'AP' resistance | $6.5\,\mathrm{K\Omega}$ |
| 'P' resistance | $2.6\,\mathrm{K\Omega}$ |
| Damping factor | 0.2 |
| Size of MTJ | $40\,\mathrm{nm}\times90\,\mathrm{nm}$ |
| Write critical current | $51\,\mu\mathrm{A}$ |

### 4.3.1 Circuit-Level Analysis

We performed the circuit-level simulation using SPICE by employing the STT model presented in [108] and TSMC 65 nm general purpose SPICE models for the CMOS implementation. Simulations are operated with a supply voltage of 1.2 V and a temperature of 27°C. The MTJ parameters used in this work are summarized in Table 7.1. Please note that our proposed techniques are generic, which can be applicable to any MTJ model.

In order to demonstrate the efficiency of our proposed technique, we simulated our circuit-level implementation at various process corners. The process variation effect is separately considered for CMOS and MTJ devices as these two are different fabrication technologies. For MTJ, $\pm3\sigma$ variations for the *Tunneling Magneto-Resistance* (TMR) and the product of resistance and area (RA) were considered. Depending upon these variation effects, the current values for fast, typical and slow CMOS process corner cases were extracted. Based on these values, the required write period for WER=$10^{-9}$ was obtained using Equation (2.13). The total time period with the segregation in terms of average delay and the necessary margin for different corner cases are shown in Figure 4.9. It is revealed from the figure that the required write margin is considerable, up to 100 % of the switching delay.

Furthermore, to study the memory architecture-level impacts, we feed the aforementioned cell-level information into NVSim tool [129]. Using NVSim, we extracted the results for various memory capacities. The results for the overall write energy and latency considering process variation are shown in Figure 4.10(a) and Figure 4.10(b), respectively. In this figure, the write latencies were directly obtained from NVSim based on various write periods specified in Figure 4.9. The energy is obtained by considering the worst-case timing scenario and variation in current due to different process corners. This is reasonable, as in a real design, the timing constraints are based on the worst-case delay. Moreover, the energy overhead due to the additional



Figure 4.9: Write period considering process variation effect for the standard, static and dynamic techniques

(a) Write latency



(b) Write dynamic energy



(c) Area overhead

Figure 4.10: Results of proposed technique with respect to various memory sizes for different process corners

circuits to generate acknowledgement signals is also considered. It is revealed from the figure that the static technique has a low write latency and energy, in comparison to the standard technique. These gains for both latency and energy are further increased with the proposed dynamic technique which is implemented on the top of the static technique. Please note that for the dynamic technique, we considered the complete timing window that is required for the step-wise current increase in order to cover the worst possible scenarios. However, in reality, some of the write operations finish earlier and in that case, the overall energy consumption would be less.

In addition to the timing and energy evaluations, we also calculated the area overhead. To estimate the area, we considered the memory architecture as shown in Figure 4.11(a). Here the bit-array is partitioned into several blocks, each of which has a dedicated periphery circuitry. Using NVSim, we obtained the ratio of the bit-array to the total area (23 % for 512 KByte). Moreover, the bit-cell layout is drawn based on the 65 nm TSMC technology rules as shown in Figure 4.11(b). The layout comparison of a single bit-cell for our proposed technique with the standard 1T1MTJ layout shows a 74 % increase in the bit-cell layout area. The addition of the PMOS transistor increases the height of the memory array (row decoder block) and the width remains unchanged (see Figure 4.11(a)). With that, the overall area overhead for static and dynamic technique is around 17 % and 20 % for a 512 KByte memory, respectively.

The area overhead of both static and dynamic techniques, after including other additional

(a) Memory architecture for write Implementation



(b) Bit-cell layout

Figure 4.11: Memory architecture for the proposed write techniques

circuits, for various memory sizes are shown in Figure 4.10(c). The area overhead of the dynamic technique also includes the area occupied by the additional circuits to generate the acknowledgement signals. According to the figure, for lower memory capacities, the area of the periphery circuits is dominating, and thus the overhead of our proposed techniques is smaller. However, with the increase in the memory capacity, the area overhead increases as well, because the bit-array area becomes more dominant. Nevertheless, the total area is far less than that of the SRAM (see Section 4.3.2).

In summary, with the proposed dynamic technique, the latency and the energy are reduced by 71 % and 43 % for a single bit-cell, respectively. The static technique has a gain in timing and energy of around 55 % and 27 %, respectively. Besides these gains, the static and dynamic techniques also impose an area overhead of less than 19 % and 22 % for a 2 MByte memory, respectively.

## 4.3.2 Comparison with state-of-the-art

In order to demonstrate the benefits of our techniques, they are compared to the column boosting technique [25], as illustrated in Table 4.3. For this purpose, we implemented this technique in our simulation framework. As shown in the table, the energy consumption for the column boosting technique for the bit-cell with the current boosting are around 781 fJ and 623 fJ for 'P'→'AP' and 'AP'→'P' switching, respectively. The energy value for the bit-cell with the current boosting was obtained by increasing the current to the level that the worst case process corner latency became equal to that of the typical process corner. However, the write current increases unnecessarily for the other bit-cells in the same column which are already fast. The energy consumption for such bit-cells (which have nominal switching latency) are increased by approximately 48 % and 33 % for 'P'→'AP' and 'AP'→'P' switching, respectively. Note that the total write period of the entire word for this technique remains the same. Moreover, the column boosting technique is a column-wise static technique which is applicable during the production test requiring a fuse based memory block organization. On the other hand, both our proposed static and dynamic techniques are bit-wise and applicable during the design

Table 4.3: Comparison of column boosting technique [25], proposed static and dynamic write techniques for a single bit-cell.

| Parameters | | Column boosting technique [25] | Static technique | Static + Dynamic technique |
|---|---|---|---|---|
| Write period [ns] | | 18.0 | 8.0 | 5.2 |
| Energy [fJ] | 'P'→'AP' | 781.1[†] and 1155.9[‡] | 436.2 | 396.3 |
| | 'AP'→'P' | 623.1[†] and 828.6[‡] | 456.0 | 357.4 |
| Technique | | Column-wise static | Bit-wise static | Bit-wise dynamic |
| Margin Addressed | | Process variation | Asymmetricity | Process variation + Asymmetricity + Stochastic |
| Applicable during | | Test phase | Design phase | Design phase |

[†]: worst process corner case bit-cell with current boosting
[‡]: typical process corner case bit-cell with current boosting

Table 4.4: Comparison of SRAM, standard STT-MRAM and proposed write technique for a 512 KByte cache

| Parameters | SRAM | Standard STT-MRAM | Static + Dynamic write technique |
|---|---|---|---|
| Area [mm2] | 2.7 | 1.31 | 1.57 |
| Read latency [ns] | 1.2 | 1.4 | 1.4 |
| Write latency [ns] | 1.2 | 19.0 | 6.9 |
| Read energy [pJ] | 574.5 | 79.0 | 79.0 |
| Write energy [pJ] | 162.3 | 1046.0 | 849.0 |
| Leakage [mW] | 966.5 | 193.0 | 193.2 |

phase. This indicates, our proposed techniques have a finer granularity, which result into the better write period and energy efficiency. On the contrary to the column boosting technique that addressed only process variation, our proposed dynamic technique additionally targets the asymmetry and stochasticity.

Furthermore, we also compared our proposed dynamic technique to SRAM for a 512 KByte cache. The comparison results for SRAM, standard STT-MRAM and our proposed techniques are shown in Table 4.4. These results are extracted using NVSim tool [129] for a 65 nm technology node with the typical process corner values. As shown in the table, the total area of our proposed technique is slightly higher than the standard STT-MRAM, but it is far less than that of the SRAM. Moreover, both standard STT-MRAM and proposed write technique have considerably low leakage compared to the SRAM as illustrated in the table. Note that, as mentioned earlier, read can be optimized independently to the write operation. Hence, we assumed the same read latency and energy for the proposed technique as the standard STT-MRAM design.

### 4.3.3 Reliability Analysis

In addition to the analysis of timing, energy and area of proposed techniques, we also performed a reliability analysis considering TDDB. The time to breakdown is very sensitive to the current flowing through the bit-cell [113, 157]. Hence, it is important to study TDDB for our proposed techniques as they rely on increased write current through the MTJ cell. To analyze the impact on TDDB, we used the model that is proposed in [114], in which this phenomenon is modeled similar to the gate dielectric breakdown in MOSFET, given as:

$$W_{hard\_breakdown} = \beta ln\ t - \beta ln\ \alpha \qquad (4.3)$$

Table 4.5: Energy reduction (%) and access latencies using our proposed technique (unit: J)

| | L1=SRAM/L2=STT-MRAM | | |
|---|---|---|---|
| Technique | Baseline | Static | Static+Dynamic |
| Latencies (L1/L2) | 1cycle/19cycles | 1cycle/9cycles | 1cycle/7cycles |
| Bzip2 | 2.45 | 2.41  (2%) | 2.40  (2%) |
| Equake | 2.75 | 2.46  (10%) | 2.43  (12%) |
| Gzip | 1.83 | 1.76  (4%) | 1.75  (4%) |
| Twolf | 1.50 | 1.25  (16%) | 1.23  (18%) |
| VPR | 3.35 | 3.09  (8%) | 3.06  (9%) |
| MCF | 1.97 | 1.71  (13%) | 1.69  (14%) |
| Average Reduction | | **9%** | **10%** |

where $\alpha$ and $\beta$ are current-dependent parameters, and $t$ is the duration of the current flow, i.e. proportional to the write period. The analysis reveals that we are able to reduce the TDDB effect by 17 % and 23 % with static and dynamic write techniques, respectively. This means that, although the write current is increased, but the reduction of the write period is more significant, and the overall impact is the TDDB improvement.

### 4.3.4 Architecture-Level Analysis

Beside the circuit-level evaluation we also performed an architecture-level analysis to obtain the influence on energy and possibly performance of an entire microprocessor. For this purpose, we employed gem5 [158] to model an embedded, in-order single-core processor with two levels of cache running at 1 GHz. The L1-Cache is split into a Data-cache and an Instruction-Cache, each of which is 32 KByte large. Our proposed technique is used only for L2-Cache, which has a size of 512 KByte, while the L1-Caches are implemented with SRAM. The corresponding access-latencies and the per-access energies were extracted with the circuit-level platform using NVSim [129]. To evaluate our techniques, we run the first 5 billion instructions of six SPEC2000 benchmarks on this system. Please note that, to reduce the overall costs of the embedded processor, the caches use a single access latency, i.e. a read access takes as much time as a write access. Our proposed technique has a significant improvement on the cache latencies, and thus overall system performance improves by 11 %. As shown in Table 4.5, our proposed technique also improves the energy reduction by 9 % and 10 % on average for the static and dynamic technique, respectively. Since the performance improvement can be used to power down the entire microprocessor more often, it also leads to significant power savings of around 11 % for the entire microprocessor.

## 4.4  Conclusions

STT-MRAM is a promising memory technology because of its beneficial features such as non-volatility, high density, scalability, high endurance and CMOS compatibility. However, the long write latency and high write energy are still major drawbacks for this memory technology. We proposed two complementary circuit-level techniques to reduce the overall write period. The first one is a static technique in which the write current is increased only for the 'AP' configuration which is the limiting factor. The second one is a dynamic technique to reduce the timing margin required to deal with the stochastic write behavior in STT-MRAM. Therefore, the write current is increased in a step-wise manner only for the bit-cells whose write operations are not completed within a given duration. With our proposed techniques, we obtained equal switching latencies for both types of transitions (3.9 ns) and we can reduce the overall write latency

by 71 %. Applying our proposed techniques to L2-cache (512 KByte) of a high-performance microprocessor improves its performance by 11 %, while its energy is reduced by 10 %.

# 5 Low-Power Multi-port Memory Architecture

In this chapter, we have utilized a fact about *Magnetic Tunnel Junction* of *Spin Orbit Torque* (SOT-MTJ) that it can inherently perform simultaneous read and write operations. In an SOT device, since read and write currents have isolated paths, both can flow at the same time through the same bit-cell without affecting each other functionality. In this way, the read-write contention can be resolved at the bit-cell level. We exploit this feature of SOT-MTJ to design a unique multi-port memory architecture in which we have modified the standard SOT bit-cell structure [159]. Overall, the simultaneous read and write feature of the SOT simplifies the multi-port design. Additionally, two low-power schemes, namely, *Unnecessary Write Termination* and *Unnecessary Write Avoidance*, are introduced on the top of our proposed multi-port memory architecture for further energy improvements. This chapter is organized as follows: related work is mentioned first and the implementation of our proposed multi-port architecture is described next. Afterwards, the comprehensive results for our proposed technique is explained, followed by the summary and conclusion is mentioned in the end of this chapter.

## 5.1 Related work

Recently, a dual port bit-cell architecture using an SOT is proposed in [160]. However, the sneaky current path during operational mode is overlooked in this architecture, which can easily result in read and write failures. Additionally, in this architecture, a level converter is also required for each write circuitry that significantly contributes to the total area and energy. Another dual-port bit-cell architecture is proposed in [161], and a multi-level cell based structure for a Graphic Processor Unit architecture is proposed in [162]. Nevertheless, both use STT memory technology that has a high switching latency and energy, due to which the overall multi-port design becomes slow and energy inefficient. Therefore, a reliable low-cost multi-port architecture which has latency and energy efficiency, is still missing.

## 5.2 Multi-port memory using SOT-MTJ

The device structure of SOT-MTJ can support simultaneous read and write operations with a negligible influence on each other. With this attribute, the read and write contention of multi-port memories can be inherently resolved at the bit-cell level and delivers a dual-port characteristic. In this section, first we explain the overview of the SOT based multi-port memory and the concept of simultaneous read-write, followed by the bit-cell architecture and the corresponding array organization. At the end, memory port enhancements and redundant write avoidance techniques for our proposed multi-port design are discussed.

### 5.2.1 Overview of multi-port memories

Multi-port memories are extensively used as shared bit-cell array for microprocessor as it can provide data access parallelism to enhance the overall performance. A block diagram of a

Figure 5.1: Block diagram of 1R1W multi-port memory using SOT bit-cell.

typical two-port memory (1R1W: one read and one write) using SOT bit-cell, is shown in Figure 5.1. It consists of two decoders to facilitate the two different addresses to access the bit-cell array and read-write periphery circuitry to perform necessary memory operations, as shown in the figure. Here read-write circuits are the same as the ones employed for the standard SOT-MRAM design. In general, our proposed multi-port design is implemented on the top of the single-port SOT-MRAM with the modification of just additional accessibility to the bit-cell. This is done since the bit-cell inherently provides the characteristics for simultaneous read and write operations, which simplifies the overall multi-port memory architecture design. In order to utilize this feature and to enhance the bit-cell accessibility, the bit-cell architecture has to be altered, which is explained in next subsection. Please note that the proposed architecture is demonstrated for 1R1W memory port, however, the number of ports can be easily extended by adding access lines and introducing corresponding bitlines in the design.

### 5.2.2 Simultaneous read-write concept

In standard single-port SOT-MRAM, the read current flows from the read terminal to any one of the write terminals through the oxide layer. In reality, this read current can pass through either of the two write terminals due to the symmetric structure of the SOT-MTJ cell. We exploit this fact about the SOT-MTJ cell, in which the read current is directed to follow the write current path in the case of simultaneous read-write operations as described in Figure 5.2(a). In other words, the conventional write current always flows from a high potential terminal to the one at a low potential, and hence, the read current can also flow through the same path after



(a) Simultaneous read-write

(b) Read and write access

Figure 5.2: Demonstration of simultaneous read and write concept on SOT-MTJ.

| Access transistors | | | Operations |
|---|---|---|---|
| N1 | N2 | N3 | |
| ON | ON | OFF | Write only |
| OFF | ON | ON | Read only |
| ON | ON | ON | Read + Write |

Figure 5.3: Modified bit-cell architecture

passing through the oxide layer. In addition, as illustrated in Figure 5.2(b), the read access is significantly faster than (around 5 X in our setup) the write access. Therefore, always old data can be read without any conflicts in the case of simultaneous read-write operations on the same bit-cell.

### 5.2.3 Bit-cell architecture modification

As mentioned earlier, due to symmetric structure of the SOT device, a read current can flow through either of the write terminals. However, in the present bit-cell structure, there is a dissimilarity in the two read current paths. One of the paths has an additional NMOS transistor (write access transistor) in the bit-cell, whereas the other path has no such transistor. Because of this, a disproportionate amount of read current flows, which can disturb the distinguishability of the two read values. This is due to the reason that the bit-cell read current value has to be compared with a reference current value which is always fixed. To avoid such scenario and to develop a symmetric read current path for better readability, we introduce an additional NMOS transistor to the other write terminal as well (connecting to the sourceline as shown in Figure 5.3). The new bit-cell architecture also serves an additional purpose that it eliminates the asymmetry in the write currents which is essentially because of the potential degradation due to the existing single NOMS transistor [154]. The addition of another NMOS transistor balances the write current, which in turn balances the write latency. This extra access transistor is activated when either (or both) read or write word line is activated. In other words, this transistor is required to be on for both read and write memory operations. Therefore, we ORed the read and write word lines, and the gate of the newly added access transistor is connected to the resulting output. The bit-cell operations with respect to their access transistors are illustrated in the table shown in Figure 5.3. The functionality of this bit-cell architecture is further explained next.

### 5.2.4 Array organization

In this section, the functionality of memory operations in term of memory array organization is explained. All the four possible scenarios for the proposed multi-port memory array accesses are described below:

**When only write is performed**

For the write operation, a bi-directional current path is established using write circuitry as illustrated in Figure 5.1. The two directions of the write current are decided based on the input data value. With the modified bit-cell architecture, the directions of the current for writing respective values are demonstrated in Figure 5.4(a). Here, we have shown the final inverters in order to illustrate the complete current source and sink paths through the bit-cell. These two inverters are the part of the write block that drives the source and write lines. During the write

(a) Write only operation

(b) Read only operation

(c) Simultaneous read-write on same bit-cell

(d) Write and read current path for simultaneous read and write on different bit-cells.

Figure 5.4: Demonstration of write and read current flow path for write only, read only and simultaneous read-write on the same cell.

operation, these two inverters have always opposite output values. In order to establish the complete current path, the two access transistors, which are connected to the write terminals, are always on during the write operation. Note that the access transistors connected to the read terminal should be off when only write operation is performed, to avoid any sneaky current path.

**When only read is performed**

For read only operations, the write access is disabled by setting write enable signal to '0'. Therefore, both inputs of the inverters of the write block become '1', which make source-line and write-line '0'. Since, only the read operation is performed, access transistor 'N3' and 'N2' are on. Therefore, a read current path is established as illustrated in Figure 5.4(b). In our implementation, we have employed a pre-charge based *sense amplifier* (SA). In this SA architecture, one end of the MTJ is connected to the amplifier component and the other end is connected to ground, so that an appropriate read current, which is evoked using SA, can flow through the MTJ device.

**Simultaneous read and write on the same bit-cell**

When both read and write operations are performed on the same cell, all three access transistors (namely, 'N1', 'N2' and 'N3') of the bit-cells are on. In such scenarios, the write current behaves same as illustrated in Figure 5.4(a). However, the read current path has to be altered based on the write current flow directions, as shown in Figure 5.4(c). In other words, the read current always follows the write current path after passing through the oxide layer of the MTJ cell. For instance, if write '1' is performed, the read current flows through the transistor 'N2'. Otherwise, it flows through the transistor 'N1' in the case of write '0'. Please note that, the sneaky write current path cannot be established in such scenarios because either source line or write line has to be at the '0' potential and the potential developed at the read terminal is always more than that. In addition, the metal electrode of the SOT device, through which the write current flows, is made up of tantalum bearing very low resistance compared to the resistance value of the MTJ device itself. Hence, the resistance value of metal electrode has negligible impact on the overall read current and readability.

In SOT-MRAM, in general, the read operation is performed much faster than write, which means there is a significant delay margin (of around 5 X in our implementation) between the read and write latencies. Therefore, during simultaneous read and write operations on the same value, the previous value in the bit-cell is always read first, before writing a new value into the bit-cell. Furthermore, the margin between read and write latency is significant enough to maintain the impact of process variation. i.e. to gurantee that read always finishes before write. In case, if there is any overlap between the slow read and fast write due to extreme influence of the process variation, the read and write circuitries can be easily tuned to make them fast and slow, respectively. Since read and write latencies are inversely proportional to their respective current values, increasing read current and decreasing write current can make read and write circuitries fast and slow, respectively.

**Simultaneous read and write on different bit-cell**

In this case, simultaneous read and write operations are performed on different bit-cells, nevertheless they share common bitlines in the same column. Bit-cells with the two independent read and write operations on the same column are illustrated in Figure 5.4(d). Similar to the read-only and write-only operations, the access transistors 'N2' and 'N3' are on for the bit-cell where read is performed. On the other hand, access transistors 'N1' and 'N2' are on for the bit-cell where write is performed. Therefore, the read current flows through the oxide layer of the MTJ cell whose value needs to read and follows the same current path as the write current of the MTJ on which the write operation is performed. The read current path with respect to the two write current path cases are shown in the figure.

### 5.2.5 Memory ports enhancement

In general, communication applications require multiple read-write ports in order to obtain high throughput in a processor. In this paper, our focus is 1R1W multi-port memory architecture, nevertheless extra read-write ports can be easily added to our implementation just by expanding the bit-cell accessibility and increasing corresponding read, source and write lines. Similar to the conventional multi-port architecture, multiple read through the same bit-cell are possible. Moreover, read and write operations on the same cell (r-w contention) can be addressed in a similar way as explained previously. Furthermore, with the addition of ports, the read currents from different ports can follow the write current in the similar way as described earlier. Otherwise, when no write operation is performed, there has to be a dedicated source-line access through which potential '0' can be provided. On the other hand, write-write contention can be resolved based on prioritizing the signals similar to the conventional methods.

### 5.2.6 Avoiding redundant writes

An SOT storing device requires a constant current value (around 130 uA in our implementation setup) to switch its magnetization, resulting in significant write access energy consumption. For an energy efficient design, we need to make sure that this architecture consumes energy only when it is necessary. Therefore, in our implementation, we have classified the write operation for each bit into the *necessary* and *unnecessary* categories, which are defined as:

- Unnecessary operations: When the value to be written in the bit-cell is already stored. In this case, the magnetic switching in SOT device does not happen.

- Necessary operations: When the value to be written in the bit-cell is different from the stored one. In this case, the magnetic switching in the SOT device actually happens.

To distinguish between the necessary and unnecessary operations, it is required to know the content of the bit-cell in which the write operation has to be performed. To do so, a read has to be performed prior every write and the write operation is performed after comparing it with the value to be written, similar to the one proposed in [133]. Unlike this technique, we require a dedicated sense amplifier and read-line in our design implementation. In this paper, we term this scheme as an *Unnecessary Write Avoidance* (UWA). With the UWA scheme, considerable energy is saved, nevertheless there is a timing penalty as the circuit delays of the read and comparator are added to the write latency. This timing penalty can be completely eliminated if the read and write operations are performed concurrently (similar to the one explained in the Section 5.2.4), which we term as an *Unnecessary Write Termination* (UWT) scheme.

## 5.3 Experimental Results

We have performed a detailed circuit-level analysis for our proposed simultaneous read-write multi-port memory architecture. In this section, we first explain the simulation setup, followed by the bit-cell analysis in which the design behavior of our proposed multi-port memory with respect to the standard single-port bit-cell design is illustrated. At the end, a multi-port memory architecture-level analysis is presented to demonstrate a comprehensive comparison of our proposed design with the conventional CMOS-based multi-port design.

### 5.3.1 Simulation setup

For the circuit-level implementation, we employed the SOT based MTJ model proposed in [163]. The MTJ parameters used for our work are depicted in Table 7.1. We used TSMC 65 nm

Table 5.1: SOT-MTJ parameters

| Parameters | Value |
|---|---|
| Damping factor | 0.5 |
| Bias magnetic field [164] | 0.1 T |
| Saturation Magnetization | $1.1 \times 10^6$ A/m |
| Metal Electrode [165] | $100\,\text{nm} \times 50\,\text{nm} \times 2.5\,\text{nm}$ |
| Critical current | 59 uA |
| 'AP'/'P' resistance | $10\,\text{K}\Omega/5\,\text{K}\Omega$ |

general purpose SPICE model for the simulation of CMOS components. Simulations were performed using Cadence Spectre tool with a power supply voltage of 1.2 V and temperature of 27° C.

### 5.3.2 Bit-cell analysis

To demonstrate the efficiency of our proposed multi-port memory architecture, we performed a bit-cell analysis using SPICE simulations. In this analysis, we have compared the behavior of the standard single-port SOT-MRAM design versus our proposed multi-port design. For the standard single-port design, we considered 2T1MTJ (two access transistor and one MTJ cell) type of bit-cell design. On the other hand, the modified 3T1MTJ (3 access transistor and 1 MTJ) type bit-cell, as described in Section 5.2.3, was used for our proposed multi-port design. Note that, similar read and write circuitries were used for both of these designs. For the proposed multi-port design, the circuit-level analysis was performed in two phases: (1) Independent operation, i.e, when only one operation (either read or write) is performed. (2) Simultaneous operation, i.e., when both read-write operations are performed at the same time.

The results for the independent and simultaneous operations for our proposed multi-port design along with the standard single-port design are summarize in Table 5.2. As shown in the table, the read latency and energy of our proposed independent operation are slightly increased compared to the standard single-port design. This is due to the addition of an extra NMOS transistor in the bit-cell of the proposed multi-port design. Note that, we have used pre-charged based sense amplifier in which only a peak of current is generated (due to a short-circuit) during a read evaluation. On the other hand, the read latency of our proposed multi-port simultaneous operations is considerably higher compared to both standard single-port and proposed multi-port independent read operations. The reason for the same is that the write current which is relatively very high compared to the read current, slows the read process. Furthermore, the write latency and energy are same for all three cases because the addition of a low read current for a very short duration has almost no impact on a high write current value. Please note that, the read and write latencies are inversely proportional to their respective current values, which means, an increase in the current value can be used for the latency reduction for both read and

Table 5.2: Proposed simultaneous read-write 3T1MTJ design versus standard single-port 2T1MTJ design

| Design parameters | | Standard single-port 2T1MTJ | Proposed 3T1MTJ | |
|---|---|---|---|---|
| | | | Independent operation | Simultaneous operation |
| Read | Latency (ps) | 28 | 30 | 46 |
| | Energy (fJ) | 0.26 | 0.34 | 31.44 |
| Write | Latency (ps) | 266 | 266 | 266 |
| | Energy (fJ) | 31.23 | 31.43 | 31.44 |
| # Cycle requires for a read and a write operation | | Two | One | |

Table 5.3: Comparison of conventional CMOS based 32X32 1R1W multi-port memory with our proposed simultaneous read-write architecture

| Parameters | CMOS 1R1W multi-port design | Proposed 1R1W multi-port design | Proposed multi-port design with Unnecessary Write Termination | | Proposed multi-port design with Unnecessary Write Avoidance | |
|---|---|---|---|---|---|---|
| | | | Nec. writes | Unnec. writes | Nec. writes | Unnec. writes |
| Read Lat. (ps) | 263 | 253 | 253 | 253 | 253 | 253 |
| Write Lat. (ps) | 241 | 431 | 431 | 431 | 539 | 539 |
| Read En. (pJ) | 1.3 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 |
| Write En. (pJ) | 1.3 | 2.5 | 2.6 | 1.8 | 2.6 | 0.8 |
| Leakage (mW) | 19.5 | 8.5 | 9.8 | 9.8 | 9.8 | 9.8 |

write operations. In our implementation, both read and write latencies are optimized for the minimum energy consumption, however, these latencies can be further improved by passing higher current, at the cost of more energy.

### 5.3.3 Multi-port memory architecture-level analysis

We have developed a $32 \times 32$ 1R1W multi-port memory for our proposed architecture together with the two proposed avoiding unnecessary write schemes presented in Section 5.2.6. Moreover, we also implemented conventional CMOS-based multi-port memory with a single-ended read and a double-ended write architecture for the same configuration. The results for the latency and energy comparisons are presented in Table 5.3. In this table, the results for *Unnecessary Write Termination* (UWT) and *Unnecessary Write Avoidance* (UWA) schemes are illustrated for both necessary and unnecessary write operations.

The read latency and energy in our proposed multi-port design are smaller compared to the conventional CMOS design. This is because, in our implementation, we employed pre-charged based SA which is fast and energy efficient. On the other hand, the write latency and energy is relatively high in our proposed multi-port design. This is due to the fact that the SOT device requires a constant current value (around 130 uA in our setup) for a specified duration to switch its magnetization. The switching latency, as mentioned earlier, can be improved with the increase in the write current value. Moreover, the energy efficiency can be significantly improved in the case of unnecessary writes using avoiding unnecessary write schemes with our proposed multi-port design. For instance, compared to proposed multi-port design, the energy in UWT and UWA can be saved for unnecessary writes by around 28 % and 68 %, respectively. The reason of such a high energy saving for UWA scheme is that the write operation is not performed at all for the case of unnecessary writes. In that case, the energy is only consumed by a read and comparator circuitry during the write operation. Nevertheless, there is a timing penalty of 25 % compared to our proposed multi-port design. Whereas, for UWT scheme, the write operation is performed along with the read and the comparison. Hence, for UWT scheme, there is no timing penalty, but the energy saving is less compared to that of UWA. In general, more than 75 % of the write operations are the unnecessary writes [132, 135]. Since our proposed multi-port design with UWT and UWA are bit-wise technique, the overall energy savings for these approaches can be significantly high at the application-level. Moreover, our proposed architecture is very effective for the leakage reduction because its bit-cell has zero-leakage unlike CMOS design. In our proposed design, the leakage is only contributed by the periphery circuitries such as decoder, read-write circuitry, drivers, etc. Moreover, for the write avoidance schemes, leakage energy is more because of the additional read and comparator circuitry.

We have also analyzed the area for the multi-port design. The proposed multi-port architec-

Figure 5.5: Word-line and bitlines metal tracks for 1R1W bit-cells of proposed multi-port bit-cell architecture.

ture requires three vertical and three horizontal metal lines and the integration of UWT and/or UWA schemes adds one extra vertical metal line as illustrated in Figure 5.5. The SOT-MTJ cell, which is small in size, is fabricated in different technology layer and its bit-cell comprises of at most four NMOS transistors for 1R1W multi-port design. Hence, the overall SOT based bit-cell is small in area, and the metal pitch of the metal lines for both rows and columns decide the total array area. We estimated the area of a $32 \times 32$ 1R1W multi-port memory using TSMC standard cell library layouts for appropriate sizes. With the addition of write avoidance scheme, due to an extra read and a comparator circuitry, the overall area is increased by 12.7 % compared to our proposed multi-port design.

In summary, our proposed SOT based multi-port memory design has less read latency and energy, and also provide significant leakage reduction. Moreover, the integration of UWT and/or UWA makes the design more energy efficient, especially for applications where switching activity is relatively low.

## 5.4 Conclusions

Spin Orbit Torque is an emerging non-volatile magnetic memory technology which has several advantageous features such as scalability, low area, high endurance, no leakage and immunity to soft-errors. In addition, this device has a capability to perform simultaneous read and write operations on the same cell without affecting each other functionality. In this paper, we propose a novel multi-port memory architecture by exploiting the simultaneous read-write attribute along with its other benefits. We demonstrated the functionality of 1R1W multi-port design architecture in all possible scenarios. Experimental results show that our proposed multi-port design is highly beneficial for read operations and leakage reduction. Furthermore, addition of avoiding redundant write schemes can increase its energy efficiency up to 68 %.

# 6 Non-Volatile Non-Shadow Flip-Flop Architecture

In this chapter, we propose a novel *Non-Volatile Non-Shadow flip-flop* (NVNS-FF), in which *Magnetic Tunnel Junction* (MTJ) cells are employed as active components for the storing/latching operations [166]. In our unique NVNS-FF design, the write process to the NV component is initiated at the same time when the input latched value is transferred to the output ports. The read value from the NV component is transferred to the output at the very next clock transition. In this way, the MTJ switching delay does not contribute to the critical path delay, and it is also effective for both short and long power-down periods. For better power saving and improved endurance, we utilize a redundant write avoidance scheme, due to which the energy efficiency in the active mode is increased to the level of conventional CMOS-based flip-flops. For the MTJ cells, we have chosen the *Spin Orbit Torque* (SOT) technology instead of the traditional Spin Transfer Torque, as it offers shorter switching delays at lower currents. In this chapter, the related work is discussed first, and then, the proposed flip-flop architecture is explained next, followed by the simulation results. In the end of the chapter, the summary and conclusion is mentioned.

## 6.1 Related work

The shrinking of device geometry and high performance requirements lead to a substantial increase in leakage current. The best way to reduce leakage current is to cut-off the power supply for certain logic blocks which are idle, is known as *power gating*. In this regard, solutions are classified into two categories: 1) volatile solutions that always require power supply for data retention, 2) MTJ based non-volatile solutions.

To perform power gating operations reliably and to speed-up the wake-up sequences, a *State Retention Power Gating* (SRPG) flip-flop can be employed [167]. In SRPG based design blocks, the master latch along with combinational logic can be power gated. However, the slave latch always require a power supply to retain the data. This, not only raises the static power consumptions, but also increases the overall area and design complexity. Another well-known technique named *Save and Restore Power Gating* (SaRPG) [168], is applied when idle periods are very long. In this technique, the flip-flop content is stored in a memory array during power-down. During the wake-up, the block is powered-up with a reset condition and the previous state is restored from the memory array. Although with this approach, the power can be completely turned-off, it incurs huge latency and area costs. Overall, CMOS based solutions are not efficient to deal with static power challenges. Therefore, researchers are looking for other alternatives such as MTJ based non-volatile solutions.

Recently, several non-volatile MTJ based flip-flop designs were introduced to reduce the static power, in which idle design blocks can be turned-off completely. Because of the non-volatility, the system state restoration is much easier than conventional CMOS flip-flops. For instance, [169] has exploited an STT-MTJ cell as a backup to store the content of the flip-flop during the power-down mode. In this design, the conventional CMOS-based flip-flop is extended with *a shadow or backup* latch which is based on MTJ devices. During power-down mode, the content of the conventional flip-flop is stored into the backup latch, and when it is power-up,

the content is loaded from the backup latch. During the active mode, only conventional flip-flop is used. Moreover, a non-volatile multi-bit flip-flop architecture [170], and several other SOT based flip-flops [164, 165, 171] have also used MTJ cells in the backup (shadow) latch. All these MTJ based flip-flop designs use MTJ as a backup storing element, which are only beneficial for the applications with a very long standby duration. Otherwise, the switching between backup and active mode adds excessive delay and energy in the design.

In summary, CMOS based solutions always require a retention power supply and MTJ based shadow flip-flop architectures provide only backup solutions which are not applicable for short power-down periods. In our proposed *Non-Volatile Non-Shadow* flip-flop architecture, we employ an SOT-MTJ cell in the active mode due to which it can be used for both short and long standby periods.

## 6.2 Non-Volatile Non-Shadow flip-flop architecture

In this section, the proposed *Non-Volatile Non-Shadow* flip-flop (*NVNS-FF*) is presented. In the beginning, we provide an overview of the flip-flop architecture and its working principle. Afterwards, we explain the basic access operations for the SOT-MTJ cells employed in our proposed flip-flop. Then, the novel NVNS-FF design is introduced followed by discussions about energy reduction, power gating and reliability.

### 6.2.1 Overview

In order to address short standby periods, an extremely fine grained power gating methodology is required, which in turn makes non-volatile flip-flops a necessity. This is achieved using our proposed NVNS-FF by employing SOT-MTJ cells as the storing elements. Our novel NVNS-FF consists mainly of two parts, i.e. a CMOS latch and a *Non-Volatile* (NV) component, as depicted in Figure 6.1(a). Similar to standard flip-flop designs, it has two output ('Q' and '$\overline{\text{Q}}$') and two input ('D' and 'CLK') pins. The main working principle of this flip-flop design, as demonstrated in Figure 6.1(b), can be divided into two phases: 1.) During the positive edge of the clock, the CMOS latch generates the outputs, and in parallel, the value stored in the CMOS latch is written to the NV component. Thus, the CMOS latch acts as a *pseudo-master* in this flip-flop architecture. 2.) During the negative edge of the clock, the write operation to the NV component is completed such that the NV component can now generate the outputs. This is used to enable the CMOS latch to capture the new incoming value from 'D'. Hence, the NV component can be seen as a *pseudo-slave*. By that means the write latency of the NV component can be hidden although the data is always stored in the NV component.



(a) Overview       (b) Working wrt. CLK

Figure 6.1: Overview of our proposed flip-flop architecture and its working principle with respect to CLK signal

Figure 6.2: Basic read write operation in SOT-MTJ

In this work, the focus is on positive edge triggered flip-flops. However, by changing the clock level, the NVNS-FF can be also used as a negative edge triggered flip-flop.

## 6.2.2 Read and write operations to the non-volatile latch

Before going into the detailed explanation of the proposed NVNS flip-flop architecture, we first explain, how the input value can be stored and read out for the NV component. This can be explained using a non-volatile latch, in which read and write operations are performed at different levels of the clock signal.

The circuit diagram for such a latch is shown in Figure 6.2. It consists of two SOT-MTJ cells, and CMOS-based read and write components. The two MTJ cells are arranged in such a way, that a self-referencing structure can be employed, which is essential for fast and reliable read operations. Therefore, the 'WT2' terminals of both cells are connected to 'CM' node. Thus, when current flows through the write terminals, it will flow from 'WT1' to 'WT2' for one cell, and in a reverse way for the other cell. By that means, always one cell stores '1', whereas the other one stores '0'.

The write (store) operation is performed using the write component of the latch, and the current directions to write a value, are illustrated in the figure. It is worth to mention that in this implementation, the write circuit component remains active only during the positive clock level. Hence, if this is the case, no write current flows through the tail transistor 'N3' because it is OFF. However, this transistor is turned-on at the negative level of 'CLK' when the read operation is performed.

The stored values in SOT-MTJ cells are read using the sense amplifier circuit as depicted in Figure 6.2. The output and its complementary values are obtained at the 'read_out' and 'read_out' nodes, respectively. Before the read is performed, these two nodes are at an equipotential using the equalizer circuit as shown in the figure. However, during the read execution, the equalizer circuit is deactivated and at the same time, the tail transistor ('N3') is turned ON. Thus, a (read) current flows through the MTJ cells. At this point, one of the output nodes instantly goes to a high steady state, while the other quickly reaches a low stable state, depending upon the resistance states of the two MTJ cells. The stabilization process is accelerated due to the two back-to-back connected inverters.

### 6.2.3 Proposed flip-flop architecture

Our proposed NVNS-FF is designed with the same read and write concept explained before for a simple latch. However, the design is modified to be clock edge sensitive, by adjusting the write part and the final output generation.

The schematic of our proposed synchronous positive edge triggered NVNS-FF is illustrated in Figure 6.3. As it can be seen in the figure, the main circuit components of the NVNS-FF are: the NV component presented in the previous subsection, an input CMOS latch and MUXes. The CMOS latch serves two purposes: 1) to provide a constant value during the write operation to the MTJ cells, and 2) to deliver the value at the NVNS-FF output ports, while it is written into the SOT-MTJ cells. As shown in the figure, the 'CLK' controlled (synchronous) data input 'D' is driven through 'IN1' and applied to the node 'DN1', and the inverted value is obtained at 'DN2'. This operation is performed during the negative level of 'CLK'. With the next transition of 'CLK', the input value is latched using 'IN2' and 'IN3'. At the same time, the corresponding value is transferred to the output ports of the NVNS-FF ('Q' and '$\overline{Q}$') and also written into the MTJ-SOT cells. Please note that these two nodes are disconnected from the input 'D' for the positive level of 'CLK'. At the next transition of 'CLK' (negative edge), the new input value is driven through the 'IN1' and the already stored value in the MTJ is read and propagated to the NVNS-FF outputs. Therefore, the read methodology for SOT-MTJ cells is the same as explained in the previous subsection. The latched value (using 'IN2' and 'IN3') and the read out value from the MTJ cells are transferred to the NVNS-FF output ports at the positive and negative levels of 'CLK', respectively. To control this mechanism, we introduced two MUXes with the select input connected to the 'CLK' signal. It is noteworthy to mention that the switching current for the two MTJs can be increased by replacing transmission gates ('T1' and 'T2') with clocked inverters.

Due to the design of the proposed NVNS-FF, the write operation to the MTJ cells has to



Figure 6.3: Circuit illustration of proposed non-volatile non-shadow flip-flop architecture.

be finished within the first half of the clock cycle. Thus, the clock period has to be longer than twice the write latency of the employed MTJ cells. However, it is important to note that this latency is not additive to the critical path. In fact, writing is performed at the same time when the values from 'DN1' and 'DN2' are propagated to the NVNS-FF outputs. Hence, the final clock period $T_{CP}$ follows Equation (1):

$$T_{CP} \geq MAX\{t_{CLK \rightarrow Q} + t_c + t_{setup}, 2t_{sot}\}, \qquad (6.1)$$

where $t_{CLK \rightarrow Q}$ is the 'CLK' to 'Q' delay, $t_c$ is the delay of the combinational logic, $t_{setup}$ is the setup time for the NVNS-FF and $t_{sot}$ is the switching delay of the SOT-MTJ cell ($<300\,\text{ps}$). As a result, in low-power embedded processors, the combinational logic delay is dominating due to the usage of slow yet low-power gates. Thus, for these processors, which typically have a clock frequency of $1\,\text{GHz}$ or less, the SOT write operation is not part of the critical path and consequently does not increase the clock period. However, if this flip-flop design is employed for high-performance processors with high clock frequencies and short combinational paths, the write latency might become critical. To alleviate this problem, the flip-flops can use an internal clocking scheme that changes the clock duty cycle, i.e. the period of the negative clock level is reduced, whereas the period of the positive clock level is increased. However, it is important to note that such high-performance processors are not the primary target of this flip-flop architecture.

The waveform for our proposed NVNS-FF is shown in Figure 6.4. We consider a clock period of $1\,\text{ns}$ and a sequence of "101" is provided to the data input 'D'. The magnetic orientations for both MTJs change at the positive edge of 'CLK' and store the opposite magnetizations. The 'read_out' and its complementary signals are generated at the negative edge of 'CLK' as demonstrated in the waveform. The final outputs 'Q' and '$\overline{\text{Q}}$' remain constant for the complete cycle and change based on the 'D' value. As shown in the waveform, the 'Q' output has a



Figure 6.4: Simulation snapshot for our proposed non-volatile non-shadow flip-flop architecture

Figure 6.5: Internal write avoidance approach for our proposed NVNS-FF

small glitch which is due to the read process. However, this glitch is so small that it cannot propagate to the next stage as illustrated by 'buffered Q'.

## 6.2.4 Avoiding redundant writes in SOT-MTJ

In our implementation, since we employ MTJ cells as active components in the design (i.e. not for backup purpose), a constant current (188 uA in our setup) is required for every flip-flop access to store the input value in the MTJ cells. However, in CMOS based flip-flop designs, only a current spike due to short circuit occurs during the transitions. Due to this reason, the NVNS-FF design has a comparatively high power consumption during active mode. To overcome this problem, we propose a circuit-level solution, in which redundant write operations are avoided. This means, if the value to be written is already stored in the MTJ cells, the write operation for the MTJ cells will not be performed. This is implemented using a new clock signal that is generated internally, as illustrated in Figure 6.5. Therefore, the 'read_out' value is compared to the value of the 'DN2' node using an XOR gate. Here, a feedback structure using two NAND gates, as shown in the figure, is employed to obtain a stable 'read_out' value. The resulting XORed value is then gated with the standard clock signal ('CLK') with the help of an AND gate, and the output signal is used as the new clock signal ('CLK_new') for the design. In other words, an internal clock gating scheme is employed, where the clock becomes active only if the input and output values of the NV component differ. In fact, not only the write operations are avoided, but also read is not performed, as the equalizer circuit remains disabled. Please note that all design components switch with the 'CLK_new' signal except the input latch components ('IN1' and 'IN3') and the MUXes. The input latch has to be operated with the original clock signal to prevent any metastability conditions in the design, and the MUXes need to use the original clock to evade any additional delay for 'CLK'→'Q' timing during the active mode when 'read_out' and 'DN2' are different.

Furthermore, it is worth to note that the principle of adjusting the positive clock pulse width mentioned in the previous subsection can not only be used for high-frequency designs to improve the critical path delay, but also to improve the energy efficiency of the flip-flop for low-frequency applications. Here, the period of the positive clock level (for a write operation) can be reduced, to minimize the time during which current is flowing through the MTJ cells, which in turn reduces the amount of energy consumed during the write operation.

## 6.2.5 Power gating methodology

Using our proposed NVNS-FF, a complete design block can be easily turned-off. Unlike the *State Retention Power Gating* (SRPG) storage, it does not require any power down controlling pin, and unlike shadow or backup architectures, the data is already stored in a non-volatile device. The power gating scheme itself can be implemented with PMOS header or NMOS footer switches. During wake-up, it is necessary that the MTJ read component is powered up before other flip-flop components as well as the combinational logic. This is required due to two reasons: 1) it boosts the wake-up sequences, and 2) it is more energy efficient due to a faster stabilizing process.

### 6.2.6 Reliability in terms of endurance

As mentioned earlier, by integrating the write avoidance scheme in our proposed NVNS-FF design, write operations are performed in the MTJ cells only when it is necessary. In other words, current is flowing through the write terminals of the SOT-MTJ cells only when a switch of the magnetization is required. Thus, this will also considerably improve the endurance of the NV component in the NVNS-FF. For instance, let us consider an embedded low-power microprocessor with a clock frequency of 1 GHz. As shown later in Table 6.4, the average switching activity for flip-flops in such designs is very low, i.e. less than 2 %. Hence, in one year the average number of write operations is less than $6.3 \cdot 10^{14}$. Since the endurance of SOT-MTJs is said to be better than the endurance of STT-MTJs ($10^{16}$ writes) [172], as no write current flows through the junction, an average lifetime of 16 years can be expected [43]. In addition, typically low-power devices are not permanently active. In fact, as laid out in [173], the active time is usually smaller than 1 %. Hence, even if a flip-flop has a switching activity of 100 %, the lifetime would be at least 30 years, and on average the lifetime is more than 100 years.

## 6.3 Experimental Setup and Results

In order to evaluate the efficiency of our proposed NVNS-FF, we performed a circuit-level analysis, followed by a comparison with the state-of-the-art. Finally, a system-level evaluation based on the circuit-level characteristics is discussed.

### 6.3.1 Circuit-level analysis

For the circuit development, we employed the setup detailed in Table 6.1 using the SOT-MTJ Verilog-A model presented in [163]. The simulations and characterization steps were performed with Cadence Spectre and Liberate, respectively. To extract the timing characteristics of the NVNS-FF, we used a slew value of 100 ps and a load capacitance of 10 fF.

Using this setup, our NVNS-FF can achieve remarkable timing values, as shown in Table 6.2. The important 'CLK'→'Q' delay is less than 100 ps, whereas a conventional CMOS flip-flop requires typically more than 100 ps. In addition, the 'CLK'→'$\overline{Q}$' is even smaller than the 'CLK'→'Q' delay, as 'Q' is obtained from a node ('DN2') that drives a higher load than that of its counter node ('DN1'). Also, the setup time values are comparable to conventional CMOS flip-flops, whereas the hold time values are slightly worse. In this regard it is important to note that for the targeted low-power systems the presented timing values are typically independent of the read and write delay to the SOT-MTJ cells, as these are usually not in the critical path. This is due to the fact that the clock frequency for these systems is low (1 GHz or less) and the

Table 6.1: Circuit-level setup

| Parameters | Value |
|---|---|
| VDD, Temperature, Process | 1.0 V, 27 °C, Typical |
| CMOS Technology | TSMC 65 nm GP |
| Damping factor | 0.5 |
| Thermal stability factor [165] | 104 |
| Bias magnetic field [164] | 0.1 T |
| Saturation Magnetization | $1.1 \times 10^6$ A/m |
| Metal Electrode [165] | 100 nm × 50 nm × 2.5 nm |
| Critical current | 59 uA |
| 'AP'/'P' resistance | 10 KΩ/5 KΩ |

Table 6.2: Timing characteristics of our NVNS-FF

| Parameters | Rise [ps] | Fall [ps] |
|---|---|---|
| 'CLK'→ 'Q' delay | 92.2 | 66.9 |
| 'CLK'→ '$\overline{Q}$' delay | 81.3 | 56.8 |
| Setup time | 38.3 | 38.3 |
| Hold time | 39.2 | 36.7 |

SOT-MTJ access is fast ($<300$ ps). Thus, *our NVNS-FF is capable of providing similar timing characteristics as a conventional CMOS flip-flop* in the active mode despite the fact of being non-volatile, i.e. the latency of the MTJ cells can be effectively hidden.

## 6.3.2 Comparison with state-of-the-art

In order to demonstrate the benefits of our proposed NVNS-FF compared to state-of-the-art flip-flops in terms of energy efficiency and area, we also implemented a conventional CMOS flip-flop and the SOT-based flip-flop proposed in [171]. The latter is a flip-flop which uses the non-volatile component only as a backup element, i.e. during active mode it works like a conventional CMOS flip-flop. The results of this comparison are summarized in Table 6.3.

The first important result is the energy consumed in active mode (including switching, short-circuit and leakage power) by each flip-flop. While the SOT-based backup flip-flop requires a similar amount of energy to the conventional CMOS flip-flop (as the components used during the active time are basically the same as in a CMOS flip-flop), our NVNS-FF consumes 4.8X more energy, when new data is captured by the flip-flop. However, as reported in Table 6.4 for various ISCAS-89 benchmark circuits and two complete microprocessors (FabScalar [174] and OpenSPARC T1), the average switching activity of the flip-flop data inputs is much lower than 30 %. This means that on average, in less than 30 % of all clock cycles new data is captured by a flip-flop. Consequently, with our scheme to avoid unnecessary write operations, the average energy consumption reduces considerably. This saving is further magnified for microprocessors running real applications, as shown in Table 6.4. For example in case of FabScalar executing SPEC2000 workloads, the average energy in the active mode consumed by all flip-flops in the design is just 5 pJ using the NVNS-FF, whereas it is 32 pJ in case of conventional CMOS flip-flops (or backup flip-flops that rely on standard CMOS architectures during active mode). Hence, our proposed NVNS-FF results in a lower average energy consumption in active mode than the other designs, if the average input switching activity is less than 8 %. In cases where the switching activity is significantly higher (e.g. S13207 in Table 6.4), a selective replacement scheme can be employed, i.e. only a subset of flip-flops is implemented using the proposed NVNS-FF. For instance, flip-flops storing data that does not require non-

Table 6.3: Comparison of various flip-flop designs with our NVNS-FF

| Parameters | Conv. CMOS | Backup SOT [171] | Proposed NVNS |
|---|---|---|---|
| Energy in active mode (fJ) | 13 | 13 | 59 and 63$^\dagger$ |
| Backup energy (fJ) | – | 66 | 52 and 56$^\dagger$ |
| Wake-up energy (fJ) | – | 34 | 5 and 9$^\dagger$ |
| Power down delay (ps) | – | 983 | 266 (worst case) |
| Wake-up delay (ps) | – | 305 | 103 |
| Transistor count | 26 | 29 | 34 and 64$^\dagger$ |
| MTJ count | 0 | 2 | 2 |

$^\dagger$: values with write avoidance.

Table 6.4: Energy consumption in active mode for conventional CMOS FF and NVNS-FFs with write avoidance for various benchmarks

| Benchmarks | Number of Flip-Flops | Avg Switching activity (in %) | Avg Energy (in pJ) | | |
|---|---|---|---|---|---|
| | | | CMOS | NVNS | Saving (in %) |
| s1423 | 74 | 6.9 | 0.35 | 0.32 | 11 |
| s5378 | 179 | 12.8 | 0.96 | 1.43 | -49 |
| s15850 | 534 | 9.3 | 2.69 | 3.10 | -15 |
| s13207 | 638 | 26.4 | 4.20 | 10.53 | -151 |
| s35932 | 1728 | 8.3 | 8.55 | 8.96 | -5 |
| Fabscalar [†] | 7563 | 1.2 | 32.58 | 5.67 | 83 |
| OpenSPARC T1 [†] | 10627 | 0.9 | 45.5 | 5.97 | 87 |

[†]: Processor functional workload

volatility (e.g. content of a branch predictor or microarchitecture flip-flops not constituting program/system state), can be implemented with conventional CMOS flip-flops. Please note that the write avoidance scheme is not feasible for standard CMOS flip-flop designs, as the main source of power consumption in active mode is short circuit during clock transitions due to the usage of transmission gates, even if the data input remains constant. Nevertheless, the power consumption in these cases is lower than in situations in which the input changes ($\approx$3X lower), which is taken into account in Table 6.4.

Another important result is the energy required for data backup to enable power gating as well as the energy required to recover the data after power gating. As our flip-flop always uses a non-volatile component to store the data, the power-down can be initiated as soon as the data is captured in the MTJs, i.e. after less than 300 ps. The required energy corresponds to the energy necessary to write into the two SOT-MTJ cells (52 fJ), i.e. the energy that is consumed during the first half of a clock cycle. In contrast, the backup flip-flop requires 3X more time to store the value in the non-volatile component and also the consumed energy is more (66 fJ). Similarly, our proposed flip-flop is also superior in terms of wake-up energy and delay. In fact, during the wake-up time the energy consumed by the NVNS-FF is more than 4X lower and basically corresponds to the energy consumed during the second half of a clock cycle in active mode. This huge advantage of our flip-flop design is due to high static currents in the backup flip-flop architecture [171], whereas the sensing scheme in our implementation has no direct current paths once the value is read out of the MTJ cells. In addition, the proposed NVNS-FF has a smaller wake-up delay, as the data does not need to be transferred from a shadow component to the main flip-flop before resuming the normal operation. Please note that the energy spent for backup and wake-up by the backup flip-flop architecture is additional to the energy consumed in active mode, whereas there is no additional energy consumption for our proposed NVNS-FF.

The last important point presented in Table 6.3 is the number of required transistors and MTJ cells, as this reflects the area required for the different flip-flop structures. This is the main advantage of the standard CMOS flip-flop as it requires only 26 transistors, compared to 29 transistors and 2 MTJs for the backup flip-flop, and at least 34 transistors as well as 2 MTJs for our proposed design. Indeed, if the scheme to avoid unnecessary write operations is employed, even 64 transistors have to be used in our design. Hence, we trade-off area to improve energy efficiency.

In summary, with the scheme to bypass redundant write accesses *our novel NVNS-FF design has a very low energy consumption which can be even lower than a conventional CMOS flip-flop in the active mode at the cost of an increased transistor count.* In addition, it *enables almost instantaneous power gating*, which further helps to reduce the overall static power consumption, as demonstrated in the next subsection.

Table 6.5: Configuration details for the experiments

| Processor | Single-core, 1 GHz, 5 Stages in-order |
|---|---|
| L1 Cache / L2 Cache | 32 KB / 512 KB Capacity |
| Execution units | ALU, CALU, FPU |
| Applications | MiBench : CJPEG and GSM |
| | SPEC2000: MCF, TWOLF and VPR |



Figure 6.6: Power gating durations for both a backup FF and NVNS-FFs

### 6.3.3 System-level evaluation

In order to illustrate the benefits of our unique NVNS-FF for power gating, we executed several system-level applications using GEM5, a cycle accurate performance simulator [158]. Therefore, we modeled a low-power embedded processor core running at 1 GHz with an in-order instruction pipeline architecture as detailed in Table 6.5. Consequently, our flip-flop design allows to power gate within one clock cycle, whereas the backup flip-flop implementation requires at least two clock cycles. However, short idle periods occur much more often than long idle periods, and as a result, a microprocessor using NVNS-FF can be power gated more effectively. As depicted in Figure 6.6, the total duration the microprocessor was power gated is on average 42 % when our proposed flip-flop is employed compared only 8 % for the scenario with backup flip-flops. Hence, *our flip-flop design allows to reduce overall leakage power by 5X*, on average.

## 6.4 Conclusions

With technology downscaling, it is very challenging to deal with static power. To reduce static power, power gating is often employed. However, many existing flip-flop architectures, such as CMOS-based flip-flops, always require an additional power supply to retain the system states, and various non-volatile backup (shadow) flip-flop designs are only efficient for long standby periods. To overcome these limitations, we proposed the novel *Non-Volatile Non-Shadow* flip-flop (NVNS-FF) design, in which a non-volatile component is employed as an active component, to enable fast power-down and power-up operations. Thus, both short and long standby periods can be efficiently addressed. In addition, the NVNS-FF can achieve similar energy and timing characteristics as conventional CMOS flip-flops, while reducing static power by 5X compared to backup NV flip-flops.

# 7 Fault Tolerant Spintronic Flip-Flop Architecture

In this chapter, a novel shadow flip-flop architecture is proposed, in which a generic *Fault Tolerant Non-Volatile Latch* (FTNV-L) is designed, to address various faults in MTJ cells [175, 176]. In this proposed FTNV-L design, several MTJ cells are structured in such a way that it can easily tolerate all single MTJ faults within a flip-flop. The remainder of this chapter is organized as follows: the proposed fault tolerant non-vlatile latch is explained next. Then, the comprehensive results about the proposed technique is described and the chapter ends with a conclusion.

## 7.1 Overview

Nowadays the static power dominates the total power consumption in a *System-on-Chip* (SoC) design, a trend that is increasing as the technology downscaling. Power Gating is the most effective methodology for the reduction of static power consumption, in which the power supply for the idle design blocks is disconnected. In this method, the conventional CMOS-based flip-flop designs, however, are not adequate as they always require a retention supply. In addition, the conventional CMOS save and restore scheme [168], where the content of the flip-flop is stored in memories during power-down and is restored during wake- up, contributes to severe delay and routing overheads. Hence, MTJ-based non-volatile flip-flop designs are becoming popular as the entire logic core can be power gated, and the data backup storage can be done locally for each flip-flop [169, 170, 177]. This way, the state restoration after power-off cycles becomes very fast and low cost. These flip-flops are very efficient for static power reduction, while maintaining the same performance for normal operation as conventional CMOS flip-flops. Nevertheless, due to several MTJ defects, as mentioned previously, the shadow latch component can either fail completely (loss of non- volatility feature) or may deliver incorrect values (loss of data integrity).

There are very effective solutions for defect and fault tolerance in array-based memories, such as main or cache memories based on built-in repair [178] and error correction mechanisms [179–181], like row/column redundancy, error detection/correction codes etc. However, these approaches are not applicable to flip-flop designs because flip-flops are placed as individual standard cells in the SoC layout. This means, failures in MTJ cells of flip-flops either lead to a very low manufacturing yield of the SoC or benefits of the leakage reduction cannot be utilized. Some techniques have arranged several non-volatile components to improve the read signal margin [182–184], nevertheless the defects are ignored. For instance, these techniques can not work if any one of the storing devices has an open defect. In [185], a robust flip-flop design is proposed to target soft-error due to radiations which only targets CMOS transistors and not MTJs. Moreover, it can not address manufacturing defects. A traditional solution to resolve this issue is *Triple Modular Redundancy* (3MR[1]) for the shadow latch. In this design, in total three shadow latch components are employed, and the output is generated based on a voting system. With this approach, although one fault per flip-flop can be tolerated, it incurs

---

[1]To distinguish between Tunnel Magneto-Resistance (TMR), we refer to Triple Modular Redundancy as 3MR

overall huge area, energy and performance costs.

## 7.2 Proposed fault tolerant non-volatile latch

This section presents our proposed fault tolerant shadow latch design using redundant MTJ cells. In Section 7.2.1, we explain the implementation of our proposed *Fault Tolerant Non-Volatile Latch* (FTNV-L) design. Later, Section 7.2.2 describes an algorithm to determine the required TMR and resistance values for our design.

### 7.2.1 Proposed FTNV-L architecture

As mentioned before, the manufacturing defects in MTJ cells are so severe that they can easily ruin the leakage benefits of the non-volatile latch, and the existing solution is not effective. Therefore, we propose a low cost solution using a novel fault tolerant MTJ-based latch design that can withstand various defects, and deliver a correct output. The implementation details of our proposed latch design, along with its functionality in the presence of all possible faults, are discussed next.

The circuit diagram for our proposed FTNV-L design is shown in Fig 7.1. It primarily consists of three components, namely, *write*, *read* and *MTJ cell arrangements*. The purpose of the write component is to store the content of the conventional CMOS flip-flop in the MTJ cells during power-down. This can be achieved by establishing a bi-directional current path such that the switching current flows through each MTJ cell. To assure the magnetic switching, the write component has to be designed in such a way that a sufficient amount of switching current for a required duration can flow through each MTJ. This current value is adjusted with the transistor widths in the write components, whereas its duration is synchronized with the 'PD_wr' period. Note that, the main requirement of this write process is that the two branches



Figure 7.1: Schematic diagram of proposed FTNV-L design

(i.e., *Branch-1 and Branch-2*) should always have a set of MTJs with opposite magnetizations. This design creates a self-referenced structure which is necessary for a proper read operation.

The read component of the design is composed of a pre-charge circuit, a pair of back-to-back connected inverters and a tail transistor. The purpose of the pre-charge circuit is to provide an equipotential at the output nodes (read_mtj and $\overline{\text{read\_mtj}}$) before the actual read is started. In our implementation, read is performed with the activation of the 'PD_rd'. During the read process, the pre-charge circuit is deactivated, and the two back-to-back connected inverters are coupled with the two branches of the MTJ sets, since the transmission gates T1 and T2 are ON. Additionally, the tail transistor 'N3' is also ON at the same time. Therefore, a current path is established, and the sensing process begins. During this sensing process, one of the output nodes goes to a low steady state, while the other remains at a high state. The two back-to-back connected inverters develop a positive feedback loop that accelerates the process of stabilizing the two output nodes.

The waveform illustration of the FTNV-L design is shown in Figure 7.2. Here, the read output, and the switching behavior of each MTJ along with its corresponding effective resistance value for each branch are shown. Read is performed at the negative level of 'CLK', where the resistance difference between two branches is important to deliver the correct output. We have obtained this waveform from real simulations, hence sometimes effective resistances have glitches when current/voltage around MTJs are changing.

The arrangement of the MTJ cells is one of the key components in our design implementation. All MTJs in each branch have the same magnetization, and as mentioned previously, the MTJs in those two branches always have the opposite magnetization. The branch in which all MTJs are in 'P' and 'AP' states are referred as *branch-P* and *branch-AP*, respectively. Each branch has a serial connection of the two parallel connected MTJs. This type of arrangement serves two purposes in FTNV-L design: (1) The parallel connection addresses short and open faults. (2) The serial connections are to increase the *ratio of the effective resistance difference* between the two branches, which we named *equivalent TMR* ($TMR_{eq}$). In other words, the flip-flop design has to meet the minimum $TMR_{eq}$ requirement during the read operation to



Figure 7.2: Waveform to demonstrate the fault-free functionility of the proposed FTNV-L design (for typical process corner, room temperature and see Table 7.1 for more setup informations).

(a) faulty MTJ of branch-AP



(b) faulty MTJ of branch-P

Figure 7.3: Branch wise demonstration of effective resistances for MTJ faults

generate the correct output. Thus, the equivalent resistance for the branch-P is given by the following equation:

$$R_{eq-P} = \frac{R_{P1} \times R_{P2}}{R_{P1} + R_{P2}} + \frac{R_{P3} \times R_{P4}}{R_{P3} + R_{P4}} \tag{7.1}$$

where $R_P$ is the resistance value of the corresponding MTJ that has 'P' magnetization. Similarly, the equivalent resistance for AP is:

$$R_{eq-AP} = \frac{R_{AP1} \times R_{AP2}}{R_{AP1} + R_{AP2}} + \frac{R_{AP3} \times R_{AP4}}{R_{AP3} + R_{AP4}} \tag{7.2}$$

where $R_{AP}$ is the resistance of the corresponding MTJ that has 'AP' magnetization. Using the above two equations, $TMR_{eq}$ is defined as:

$$TMR_{eq}(\%) = \frac{R_{eq-AP} - R_{eq-P}}{R_{eq-P}} \times 100 \tag{7.3}$$

If one MTJ cell has a permanent or temporal defect, the equivalent resistance changes based on the fault type, as discussed next.

**Short fault**: When one of the MTJs has a short fault, a relatively high current flows though that defective MTJ. Consequently, the MTJ which is in parallel to the shorted one, is bypassed for both read and write operations. Hence, the equivalent resistance for both 'P' and 'AP' is :

$$R_{eq-short\{P,AP\}} = \frac{R_{eq\{P,AP\}}}{2} \tag{7.4}$$

where $R_{eq\{P,AP\}}$ is the equivalent resistance of either branch-P or branch-AP.

**Open fault**: When one of the MTJs is open, no current flows through that MTJ. Unlike for shorts, the MTJ which is in parallel to the defective MTJ is usable and it becomes in series with the other two parallel connected MTJs. In this case, the equivalent resistance for both 'P' and 'AP' configuration is:

$$R_{eq-open\{P,AP\}} = \frac{R_{eq\{P,AP\}}}{2} + R_{\{P,AP\}} \tag{7.5}$$

where $R_{P,AP}$ is the resistance of a single MTJ in either 'P' or 'AP'.

**Stuck-at-P**: When one of the MTJ cells is stuck at the 'P' configuration, then only the 'AP' branch is affected. Therefore, the equivalent resistance in 'P' branch is as follows:

$$R_{eq-stuck-at-P} = \frac{R_{AP}}{2} + \frac{R_P \times R_{AP}}{R_P + R_{AP}} \tag{7.6}$$

where $R_P$ and $R_{AP}$ are the resistances of the MTJs when they are in 'P' and 'AP' configurations, respectively. **Stuck-at-AP fault**: When one of the MTJ cells is stuck at the 'AP' configuration, then only the 'P' branch is affected. Therefore, the equivalent resistance in 'AP' branch is as follows:

$$R_{eq-stuck-at-AP} = \frac{R_P}{2} + \frac{R_{AP} \times R_P}{R_{AP} + R_P} \tag{7.7}$$

All aforementioned faults with their effective resistances are demonstrated in Figure 7.3. When the faulty MTJ is in the branch-AP, the short fault has the worst effective resistance, which in-turn results in the worst $TMR_{eq}$ for that fault. On the other hand, when the faulty MTJ is in the branch-P, the open fault becomes critical from $TMR_{eq}$ point of view, compared to all other faults. For more details, please see Section 7.3.2.

## 7.2.2 Algorithm to obtain TMR and resistance values

The TMR and resistance values for any MTJ device have to be fixed beforehand at device-level prior to the circuit design implementations. In general, high TMR is always preferable for the read, however, it has some limitations due to materials, and trade-off with device parameters

---

**Algorithm 1** Algorithm to determine TMR and resistance values for MTJ device

**Input** $R_{min}$, $R_{max}$, $TMR_{min}$, $TMR_{max}$, $R_{array}$,
    $TMR_{eq\_acceptable}$, $R_{step}$, $TMR_{step}$
**Result** $TMR_{MTJ}$, $R_{MTJ}$;
**For** $TMR_{MTJ}$ **from** $TMR_{min}$ **to** $TMR_{max}$ **in** $TMR_{step}$ **do**
  **For** $R_{MTJ}$ **from** $R_{min}$ **to** $R_{max}$ **in** $R_{step}$ **do**
    $R_P = R_{MTJ}$;
    $R_{AP} = R_{MTJ} * (1 + TMR_{MTJ})$;
    Obtain the equivalent resistances for all four fault cases
using Eq (7.4), (7.5), (7.6), and (7.7);
    Obtain the $TMR_{eq}$ in each cases;
    **if** ($TMR_{eq}$ of {Open && Short && Stuck-at-AP && Stuck-at-P)}
    $\geq TMR_{eq\_acceptable}$) **do**
      **return** $TMR_{MTJ}$;
      $R_{array}$ collects corresponding resistance values ;
    **done**
  **done**
**done**
**return** $R_{MTJ}$ from $R_{array}$ that supports MTJ design parameters;

---

such as switching energy and thermal stability. Therefore, a methodology is needed to obtain a minimum TMR of a single MTJ with a design-suited resistance value, which can tolerate all aforementioned MTJ faults in the design.

To obtain the TMR and resistance values for each MTJ, a generic algorithm is developed, shown in Algorithm 1. As inputs, we consider the range of TMR and resistance values that can be supported at technology level. In general, the TMR value can easily reach more than 600 % depending on the oxide layer thickness and area of the device [20]. In addition, the minimal acceptable $TMR_{eq}$, which can generate the correct output during read, is also part of the input parameters. Here, TMR and resistance values are varied by a specific step size. For each TMR and resistance values, the effective resistance and $TMR_{eq}$ are obtained for every fault type using Equations (7.4)-(7.7). If the obtained $TMR_{eq}$ is equal to or more than the acceptable $TMR_{eq}$ value, we store the corresponding resistance values in an array, so that an optimum value, based on the device trade-offs, can be picked. Note that the range of the input TMR values can be further increased by adding another set of parallel MTJs for both branches in the design. In that case, the read and write components have to be designed accordingly. Moreover, if the write drivers are not able to provide sufficient current in such cases, it is also possible to add multiple drivers at an intermediate stage of the design. Another possibility is to use a high supply voltage to pass a high switching current.

## 7.3 Experimental Setup and Results

We performed a circuit-level analysis in order to evaluate the efficiency of our proposed FTNV-L design. The simulation setup is discussed first, followed by the circuit-level results. In the end, a comparison of our proposed technique with triple modular redundancy is performed.

### 7.3.1 Simulation setup

For the circuit design implementation, we employed the MTJ model presented in [186], and the other design parameters for the simulations are depicted in Table 7.1. Here, our MTJ model is tuned for the TMR and resistance values specified in the table, which are determined using Eq (7.4)-(7.7), with the assumption of an acceptable $TMR_{eq}$ of 50 %. We have used Cadence Spectre tool for circuit simulations.

The resistance value associated with each MTJ is obtained by measuring the current value and voltage across its terminals. Furthermore, to obtain a setup for the defective MTJ cell, we employed a resistance device to replace the MTJ in the design. For instance, a low (around 5 Ω) and a high (around 5 MΩ) resistances are connected to demonstrate the short and open faults, receptively. Similarly, to show the stuck-at-P and stuck-at-AP behavior in the design, a resistance value equivalent of $R_P$ and $R_{AP}$ is connected, respectively. Please note that only one resistance at a time is connected, as our design targets a single fault per latch. For

Table 7.1: Circuit-level setup

| Parameters | Value |
|---|---|
| VDD and Temperature | 1.2 V and 27 °C |
| CMOS Technology | TSMC 65 nm GP |
| Thermal stability factor | 60 |
| Free/Oxide layer thickness | 1.84/1.48 nm |
| RA | 6.145 $\Omega\mu m^2$ |
| TMR @ 0 V | 200 % |
| 'AP'/'P' resistance | 3.6 KΩ/1.2 KΩ |

Figure 7.4: TMR values in the presence of various faults



Figure 7.5: Read latency values in the presence of various faults

process variation, we have considered MTJ and CMOS components separately as these two are different fabrication technologies. For MTJ components, we used statistical monte-carlo model that includes variation in terms of TMR and the product of *Resistance and Area* (RA). On the other hand, for CMOS components, we used the statistical model provided by the TSMC.

### 7.3.2 Circuit functionality analysis

The TMR value is very sensitive to MTJ defects, which in turn influences the functionality of the design. In our proposed design, we have performed a detailed $TMR_{eq}$ analysis and the results for both branch-P and branch-AP are demonstrated in Figure 7.4. These $TMR_{eq}$ values influence the delay of the restore (read latency) operation of the FTNV-L as demonstrated in Figure 7.5. As shown in the figure, the worst $TMR_{eq}$ value is obtained for open and short faults in branch-P and branch-AP, respectively. The range of process corner variations in our $TMR_{eq}$ and read latency results are shown with error bars. The functionality of the FTNV-L design in the presence of the short-AP and open-P faults, are demonstrated in Figure 7.6. In these figures, the read outputs and effective resistances for both branches are shown. Here, both of those figures show the low TMR values (marked by the blue circles) during the read operation. The low TMR value is for the low resistance and high resistance value range for the short-AP and open-P faults as described by Figure 7.6(a) and Figure 7.6(b), respectively. The reason for low $TMR_{eq}$ for short-AP fault is that the effective resistance of branch-AP becomes low, close to that of the branch-P, due to an MTJ short. Similarly, in case of an open-P fault, the effective resistance of branch-P becomes high, close to that of the branch-AP, due to an

(a) Short-AP fault



(b) Open-P fault

Figure 7.6: Functionality of FTNV-L in the presence of short-AP and open-AP faults (for typical process corner, temperature of 27°C and see Table 7.1 for setup information). Blue dotted circle indicates the worst resistance differences during read for corresponding fault cases.

MTJ open. However, these low TMR values are still good enough to read the output correctly.

On the other hand, in the presence of short-P and open-AP faults, the $TMR_{eq}$ value is high, even more than that of a fault-free MTJ cell. This is because, the faulty MTJs in these two cases are additive to the resistance differences which further increases the overall effective resistance. Moreover, the $TMR_{eq}$ value for both stuck-at faults is slightly less than that of the TMR value of an MTJ. Since, the read latency is inversely proportional to the $TMR_{eq}$ value, short-P has the lowest and short-AP has the highest delay.

For write operation, the voltage drop due to the series-parallel connections of MTJs for our proposed design is similar to the standard latch design. This is because, the overall effective resistance between the two write drivers are same (see Equations (7.1) and (7.2)). However, in our proposed design, the write current is divided into two branches because of the parallel connections of the MTJs. To compensate this, the write drivers of our proposed design are strengthened (3.4 X) to deliver more switching current (around 2 X) compared to the standard design. In our proposed design, the write drivers are delivering high enough current in the

(a) Short-P

(b) Short-AP

(c) Open-P

(d) Open-AP

(e) Stuck-at-P

(f) Stuck-at-AP

Figure 7.7: Impact of process variation (1000 monte-carlo simulations) on the effective TMR values for each fault.

presence of any faults to ensure the necessary switching.

### 7.3.3 Process variation analysis

Similar to CMOS fabrication, the MTJ cell manufacturing and measurement processes also exhibit variations. In other words, due to manufacturing process, the MTJ critical dimensions such as surface area, oxide thickness, size of the free layer etc., are not the same. In general, the read latency of a non-volatile latch is influenced by process variation in two ways: 1) The resistance of the MTJ cell varies with process variation, affecting the read current that in turn the read latency. 2) Process variation affects the resistance difference (TMR), resulting in variation in the read latency, i.e., higher the TMR, lower the read latency and vice-versa. Moreover, also CMOS variation affects read by changing the read current mostly due to the transistor threshold variations. In order to perform a process variation analysis for our proposed FTNV-L architecture, we have run monte-carlo simulations for the effective TMR value ($TMR_{eq}$). The histogram results for all faults for 1000 samples are depicted in Figure 7.7. As illustrated in the figure, $TMR_{eq}$ variations for all fault cases show normal distribution. The short-P and

open-AP faults can have a very high overall $TMR_{eq}$ value as described earlier. For instance, the $TMR_{eq}$ value for short-P fault can reach upto $586\%$. However, these two fault cases have much wider distributions as shown in Figure 7.7(a) and Figure 7.7(d). Here, the $\sigma$ value for the $TMR_{eq}$ value for these two faults is more than 20. On the contrary, all other fault cases, have relatively narrow distribution with $\sigma$ value less than 10. The short-AP fault, which has the lowest TMR value among all faults, has lowest $\sigma$ value (less than 7). The reason for this behavior is that the effective resistances in both branches in short-AP are varying mostly in the same direction. In other words, the difference in the variation of the two set of resistances of the two branches is less (around $2\,X$), resulting in low TMR variations. Nevertheless, this difference is significant for short-P where the effective high resistance value can vary more than $5\,X$ compared to that of the low resistance branch value. Despite the low TMR variations for short-AP fault, the $TMR_{eq}$ value can go as low as $29\%$. With this TMR value, as explained earlier, the data integrity can be maintained, but the read latency is increased significantly (can reach up to $105\,ps$).

### 7.3.4 Variations due to different operating temperature

The TMR value of an MTJ device varies with the operating temperature as the MTJ resistances are sensitive to temperature. Therefore, in order to evaluate the efficiency of our proposed FTNV-L design for various operating temperatures, we have performed TMR and read latency analysis for all faults at different temperatures. The results are extracted for both branch-P and branch-AP for a range of temperature values as illustrated in Figure 7.8. As shown in Figure 7.8(a) and Figure 7.8(b), the TMR value decreases in each fault case with the increase



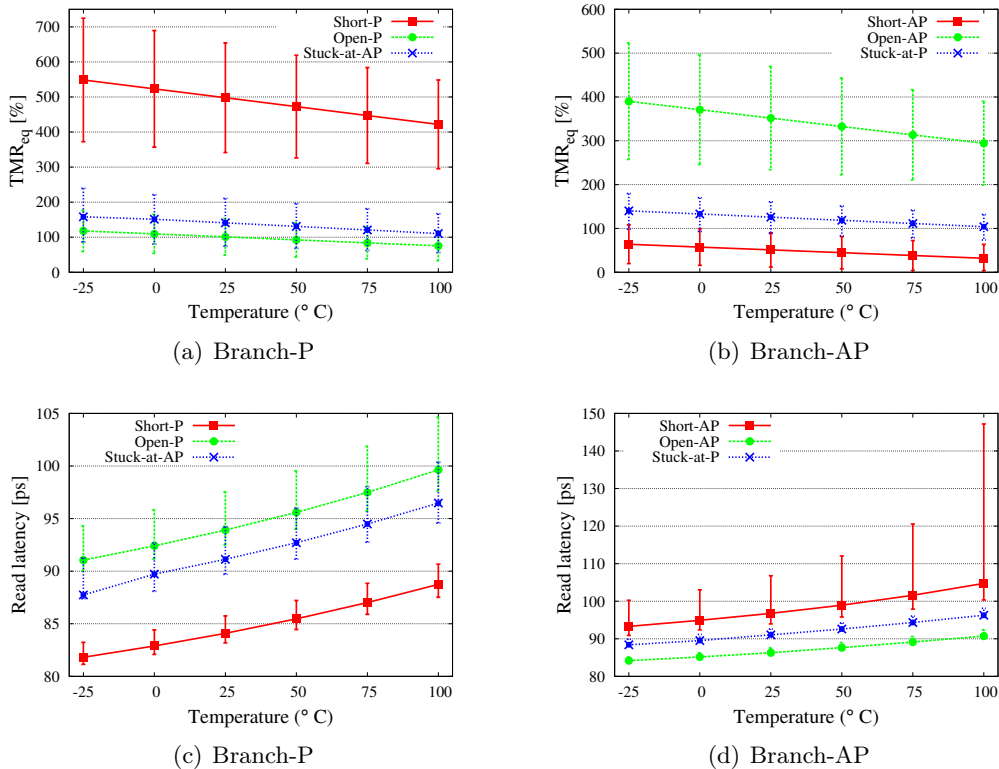(a) Branch-P

(b) Branch-AP

(c) Branch-P

(d) Branch-AP

Figure 7.8: TMR and read latency variations for various MTJ faults with respect to operating temperature.

in temperature. This rate of TMR declination with the increase in temperature is high for short-P and open-AP compared to the other faults. In addition, the range of TMR variations due to process variation is also high for these two fault cases. Nevertheless, the $TMR_{eq}$ value remains significantly high at $100\,^{\circ}\mathrm{C}$. On the other hand, the open-P and short-AP fault, which have the worst case TMR values for 'branch-P' and 'branch-AP', respectively, demonstrate relatively less TMR declination rate and variations with the increase in temperature value. Furthermore, the TMR value of the short-AP fault can reach as low as less than $5\,\%$ at $100\,^{\circ}\mathrm{C}$. The TMR variations for each fault case influence the read latency as illustrated by Figure 7.8(c) and Figure 7.8(d). The short-P and open-AP, which have best TMR values, resulting in a very low read latency (less than $93\,\mathrm{ps}$). On the contrary, the short-AP fault, which have the worst TMR value among all faults at $100\,^{\circ}\mathrm{C}$, can have very high read latency (upto $150\,\mathrm{ps}$).

Overall, using our proposed FTNV-L architecture, the flip-flop functionality remains intact, since it is able to deliver good enough $TMR_{eq}$ value for all fault cases in the presence of process variation at different operating temperatures. In some cases, such as short-AP or open-P, the $TMR_{eq}$ value found to be considerably low resulting in a high wakeup delay. However, this increase in wakeup delay is very low compared to the power down durations. Hence, it has overall negligible impact on the overall performance of the system. Beside this, in case, the two read states are not distinguishable due to extreme variation conditions, the TMR value of individual MTJs can be increased (see Algorithm 1). If the TMR value of individual MTJs has already reached to its maximum due to the material or write current limitations, then another set of two parallel MTJs can be connected to the given structure to further increase the effective TMR value.

### 7.3.5  Area analysis

In addition to the design parameters, we conducted an area analysis for our proposed FTNV-L design. Compared to the standard latch design, the only components added in our proposed design is the replacement of the two MTJs with the eight parallel/serial connected MTJs (see Figure 7.1). Nevertheless, in general, MTJs are fabricated in another layer [187], and additionally, flip-flops are widely distributed all over the logic core unlike memory bit-cells. Therefore, there would be no placement restrictions for MTJs for flip-flop designs as those can be easily placed above CMOS device layers, as illustrated in Figure 7.9. The area of CMOS layers which also includes the conventional flip-flop design (i.e., X * Y) is significantly more than the area of the magnetic layer (i.e., X´ * Y´). Therefore, the area of the CMOS layer eventually contributes to the total area of the flip-flop design. In case the magnetic layer area is more than the CMOS part of the flip-flop, e.g. the MTJs are placed in more relaxed manner or any other manufacturing constraints, they can be placed above the neighboring combinational logic cells as well, i.e., without impacting the chip area. Please note that the MTJ via sizes are negligible compared to the area of the CMOS layer. Due to this fact, the effective area of our proposed FTNV-L design remains the same as the standard latch design. However, the parallel/serial structure of MTJs in our proposed design requires an increased write current (around $2\,\mathrm{X}$) to ensure the switching in each MTJ. Therefore, the drive strength of the write component has to be increased by $3.4\,\mathrm{X}$ compared to the standard latch design, which is, however, negligible in a custom layout design for a flip-flop.

### 7.3.6  Comparison with triple modular redundancy

To illustrate the advantages of our proposed FTNV-L design, we compare it with a standard latch as well as 3MR. For the standard NV latch implementation, we use only two MTJs, one per each branch. For the 3MR implementation, we employ three standard NV latch designs
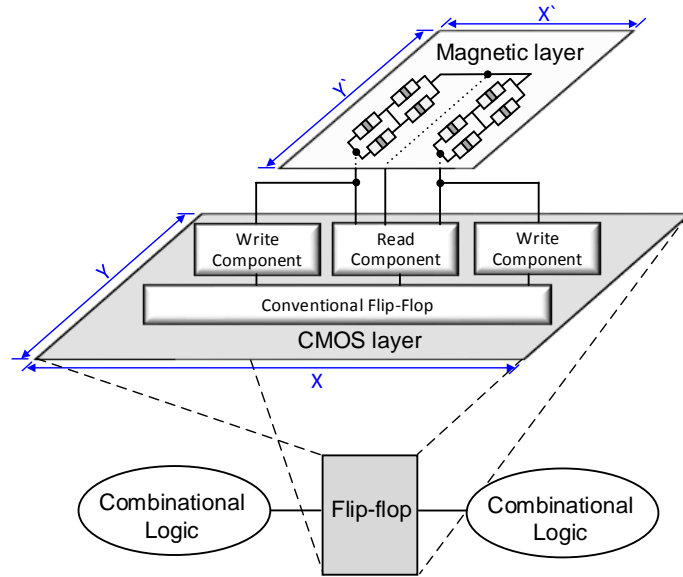
Figure 7.9: Demonstration of magnetic and CMOS layers for area analysis (area of magnetic layer $<$ area of CMOS layer).

with a voting circuit. The results of comparison for the three designs for a normal operation are summarized in Table 7.2.

As specified in Table 7.2, the $TMR_{eq}$ value for each design is the same for the normal operation when fault-free MTJs are considered. However, in the presence of defective MTJs, the standard latch design is not functional at all, whereas our proposed FTNV-L design as well as 3MR are able to generate fault-free output. Both of these designs can address a single MTJ fault per latch, but the 3MR has huge overheads because it uses three sets of standard NV latch designs and a voter circuit. For instance, compared to FTNV-L, the 3MR design has around 2 X and 3 X overheads for the read latency and energy, respectively. Here, the read operation means the MTJ values are sensed in an NV shadow latch (during wakeup) and restored back to the normal flip-flop. The voter circuit in 3MR adds 120 ps to the delay and consumes 6.4 fJ energy during restore. In the store operation, the flip-flop data is written into the NV shadow latch during the power-down mode. The storing delay is same for those three designs because the three sets of MTJs in 3MR design can be written in parallel. But, similar to read, it has around 3X more energy as it requires total three sets of MTJ to switch its magnetization.

On the other hand, our proposed FTNV-L has almost similar results in comparison to the standard NV latch design. For instance, the backup energy for a non-volatile latch is dominated by the write current as a continuous constant current requires to flow through the MTJs for a certain duration to switch their magnetizations. The equivalent resistance of our proposed design is also high because of the serial-parallel MTJ connections. To be more precise, the equivalent resistance of each branch for our proposed design becomes twice as that of the standard design (see Equation 7.1 and Equation 7.2). Therefore, we need to pass an increased current for our proposed design in order to equalize the switching latency to that of the standard latch design. To achieve that, the drive strength of the write components are increased by 3.4 X in FTNV-L design compared to the standard NV latch design. Since similar write currents flow in both designs, the write energies are comparable. On the other hand, we are dealing at the processional switching regime where the switching current can vary significantly for almost the same latency. Due to this, in automated simulated environment, with our proposed design, we could able to attain the same latency with slightly lower current value, hence we have lower total energy compared to that of the standard latch design. Moreover, the increase in the width

Table 7.2: Comparison of standard latch and 3MR design with proposed FTNV-L design

| Parameters | Standard NV Latch | 3MR | Proposed NV Latch |
|:---:|:---:|:---:|:---:|
| $TMR_{eq}$ (%) | 200 | 200 | 200 |
| Read Latency (ps) | 83 | 203 | 89 |
| Read Energy (fJ) | 12 | 42 | 15 |
| Store delay (ps) | 4056 | 4056 | 4065 |
| Store energy (fJ) | 390 | 1170 | 366 |
| Leakage (nW) | 37 | 155 | 97 |
| Transistor count | 16 | 72 | 16 [†] |
| MTJ count | 2 | 6 | 8 |

[†]:Transistor width are increased by 3.4 X compared to standard NV-latch

of transistors in the write component results in high leakage compared to the standard NV latch design as illustrated in the table. Please note that, the leakage due to the read component for our proposed design remains same as the standard latch design. Similar to previous cases, the 3MR design has a high leakage due to more circuitries as described before.

## 7.4 Conclusions

Nowadays, spintronic based shadow latches are gaining attention as these are highly beneficial for leakage reduction. This is because, the storing devices of these latches, which are *Magnetic Tunnel Junction* (MTJ) cells, have attractive attributes such as zero leakage as well as high access speed. Consequently, these latches very effective for instant-on/normally-off computing in an SoC. However, these MTJ cells are highly susceptible to several manufacturing defects such as short oxide, open vias, magnetic orientation of the free layer does not switch, etc. Due to this, the yield of the design is affected since a single defect in a flip-flop can lead the failure of the entire backup strategy. Therefore, we proposed a *Fault Tolerant Non-Volatile Latch* (FTNV-L), in which MTJs are serially and parallelly connected in a unique way to tolerate MTJ related faults. We have demonstrated the functionality of our proposed design in the presence of all MTJ faults under the influence of process variation and operating temperature. In addition, using the FTNV-L design, any single fault per latch can be tolerated at much reduced costs compared to the traditional solution based on triple modular redundancy.

# 8 Conclusions and Outlook

With the increasing performance requirements in a microprocessor, the demand of memories is continuously growing. However, the conventional memories such as SRAM, DRAM and Flash are facing severe scaling challenges with the shrinking technology nodes. As per *International Technology Roadmap for Semiconductors* report, the storing devices for these conventional memory technologies can not be further scaled beyond the year 2025 [40]. This is due to the fact that, with smaller devices, the increasing effect of process variation severely impacts several memory design constrains such as read, write and hold margins, as well as exaggerates the failure rates due to external noises and radiations. On the top of that, leakage power, which becomes a dominating factor in the total power consumption, also increases with the technology downscaling. This is due to the fact that SRAM-based memories such as caches, registers, flip-flops, latches, always require a constant supply voltage to retain their content. Moreover, for DRAM, the periodic refresh rate is increasing significantly with the down-scaling, that can further add to the total power consumption even when it is not performing any operation.

To address the challenges associated with conventional memory technology, researchers both in academia and industry have been seeking for non-volatile memory solutions. *Magnetic Random Access Memories* (MRAM) technology such as *Spin Transfer Torque* (STT) or *Spin Orbit Torque* (SOT), is a promising candidate because of its various beneficial features such as high density, high endurance, non-volatility, CMOS compatibility and immunity to radiations. The MRAM is the spintronics storing technology, in which the spin of the electron is exploited to store the information unlike the conventional memories where the charge of the electron is used for the same purpose. Overall, this memory technology combines the attributes of speed and endurance of SRAM, high density of DRAM and non-volatility of Flash. Hence, it has potential to become a universal memory technology that can be integrated into every level of the memory hierarchy of the computing system.

Despite all these merits, spintronic-based memories have several challenges which need to be addressed before widespread commercial utilization. For instance, the write access energy and latency are relatively higher than CMOS-based memory technology. This is because, it requires a significant timing margins to account for its stochastic switching as well as the impact of process variation. Moreover, because of such a high timing margin, read and write current flow for longer durations, which can easily lead to read disturb and other degradation issues such as *Time Dependent Dielectric Breakdown* (TDDB). Additionally, reliable high performance and energy efficient solutions for registers and flip-flops are still a major challenge for this technology.

In this thesis, several design techniques were proposed to address the aforementioned challenges associated with the spintronic-based caches, registers and flip-flops. Here, the fundamental properties of the storing cell were exploited using circuit design first and those were evaluated using a system-level platform afterwards. In order to make this technology energy efficient, self-timed read and write techniques were proposed, in which the respective operation completions are detected dynamically. Using these techniques, a significant amount of energy was saved as it is consumed only during when bit-cells are actually performing operations. Moreover, the impact of TDDB and read disturb can be reduced significantly using these techniques because the duration of the current flow is shortened. On the top of these techniques, a write speed acceleration technique was proposed, in order to improve the overall write period of

the memory. Using this technique, the write current is boosted in a controlled way only during the required occasions and only for bits for which it is necessary. For instance, the bit-cell asymmetry is addressed by boosting current statically only for the slow-writes. Whereas, the write margin due to stochastic switching and process variation is reduced by increasing write current dynamically in a step wise manner only for bits which have unfinished transitions. Overall, using combined self-timed and write speed acceleration techniques, the cache performance and energy efficiency were significantly improved, and also their reliability is enhanced extensively.

Besides caches, solutions for registers and flip-flops were also covered in this thesis work. For registers, a novel multi-port memory architecture was proposed. This technique was designed by exploiting a unique property of the SOT cell that it can perform both read and write operations at the same time. In this way, the read-write contention can be naturally resolved, that resulted in a much simplified multi-processing design. For flip-flops, a non-volatile spintronic-based flip-flop design was proposed which can greatly save static energy consumption. This is due to the fact that, this design can address short standby durations along with long periods of inactivity, unlike the conventional shadow flip-flop architectures. Moreover, this architecture achieved similar timing characteristics as CMOS-based flip-flops. Additionally, we have also developed a fault tolerant non-volatile flip-flop architecture, which is resilient to various manufacturing faults. In this technique, several storing cells are arranged in a unique way that can deliver a fault-free functionality in the presence of different faults. This flip-flop design is applicable to both shadow as well as non-shadow architectures.

Overall, in this thesis, significant improvements for performance, energy and reliability were demonstrated for spintronic on-chip storages such as caches, registers, flip-flops and latches. This spintronic technology can be further explored in many aspects. For instance, this technology can be developed to store multi-bits in order to enhance the density. A few techniques have been proposed to implement the spintronic-based storage cell with multi-bit capabilities, but the fundamental challenges related to the structure of the cell are yet to be resolved for the real product development. Moreover, spintronic technology can also be very beneficial in the security field. In this field the variations due process and stochastic switching can be easily utilized for the development of *Physically Unclonable Function* and *True Random Number Generator*. Moreover, the stochastic switching nature of spintronic technology can be exploited for the development of stochastic-memristive synapses for neuromorphic computing paradigms. Additionally, this technology can play a vital role for various sensor applications such as probes for biomedical purposes, magneto-resistive sensors for non-destructive testing, power measurements, fringe-fields detection for bio-chips, compass applications especially for automotive, etc.

# Bibliography

[1] "Samsung. [Online]." Available: http://www.samsung.com/us/aboutsamsung/our_company/history/timeline/, [accessed November 2016].

[2] "Hynix. [Online]." Available: https://www.skhynix.com/eng/about/history2000.jsp, [accessed November 2016].

[3] K. Itoh, M. Horiguchi, and H. Tanaka, *Ultra-Low Voltage Nanoscale Memories (Series on Integrated Circuits and Systems)*, 2007.

[4] "Intel. [Online]." Available: http://ark.intel.com/, [accessed November 2016].

[5] "As Nodes Advance, So Must Power Analysis [Online]." Available: http://semiengineering.com/as-nodes-advance-so-must-power-analysis/, [accessed February 2017].

[6] N. H. Weste and D. Harris, *CMOS VLSI Design: A circuit and System Perspective*, 3rd ed.

[7] M. N. Baibich, J. M. Broto, A. Fert, F. N. Van Dau, F. Petroff, P. Etienne, G. Creuzet, A. Friederich, and J. Chazelas, "Giant magnetoresistance of (001) Fe/(001) Cr magnetic superlattices," *Physical review letters*, vol. 61, no. 21, p. 2472, 1988.

[8] K. S. Yoon and J. P. Hong, "Observation of giant magnetoresistance in CoFeN/AlOx/CoFeN magnetic tunneling junctions employing a nitrogen-doped amorphous CoFeN free layer electrode," *Applied Physics Letters*, vol. 110, no. 1, p. 012402, 2017.

[9] J. S. Moodera, J. Nowak, and R. J. van de Veerdonk, "Interface magnetism and spin wave scattering in ferromagnet-insulator-ferromagnet tunnel junctions," *Physical Review Letters*, vol. 80, no. 13, p. 2941, 1998.

[10] H. Wei, Q. Qin, M. Ma, R. Sharif, and X. Han, "80% tunneling magnetoresistance at room temperature for thin Al–O barrier magnetic tunnel junction with CoFeB as free and reference layers," *Journal of applied physics*, vol. 101, no. 9, p. 09B501, 2007.

[11] D. Wang, C. Nordman, J. M. Daughton, Z. Qian, and J. Fink, "70% TMR at room temperature for SDT sandwich junctions with CoFeB as free and reference layers," *IEEE Transactions on Magnetics*, vol. 40, no. 4, pp. 2269–2271, 2004.

[12] H. Kano, K. Bessho, Y. Higo, K. Ohba, M. Hashimoto, T. Mizuguchi, and M. Hosomi, "MRAM with improved magnetic tunnel junction material," in *Magnetics Conference, 2002. INTERMAG Europe 2002. Digest of Technical Papers. 2002 IEEE International*, p. BB4, 2002.

[13] S. Parkin, "Magnetic Tunneling Junction Devices for non-volatile random access memory," in *Magnetics Conference, 1999. Digest of INTERMAG 99. 1999 IEEE International*, pp. GA01–GA01, 1999.

[14] M. Bowen, V. Cros, F. Petroff, A. Fert, C. Martınez Boubeta, J. L. Costa-Krämer, J. V. Anguita, A. Cebollada, F. Briones, J. De Teresa *et al.*, "Large magnetoresistance in Fe/MgO/FeCo (001) epitaxial tunnel junctions on GaAs (001)," *Applied Physics Letters*, vol. 79, no. 11, pp. 1655–1657, 2001.

[15] S. Yuasa, A. Fukushima, T. Nagahama, K. Ando, and Y. Suzuki, "High tunnel magnetoresistance at room temperature in fully epitaxial Fe/MgO/Fe tunnel junctions due to coherent spin-polarized tunneling," *Japanese Journal of Applied Physics*, vol. 43, no. 4B, p. L588, 2004.

[16] S. Yuasa, T. Nagahama, A. Fukushima, Y. Suzuki, and K. Ando, "Giant room-temperature magnetoresistance in single-crystal Fe/MgO/Fe magnetic tunnel junctions," *Nature materials*, vol. 3, no. 12, pp. 868–871, 2004.

[17] S. S. Parkin, C. Kaiser, A. Panchula, P. M. Rice, B. Hughes, M. Samant, and S.-H. Yang, "Giant tunnelling magnetoresistance at room temperature with MgO (100) tunnel barriers," *Nature materials*, vol. 3, no. 12, pp. 862–867, 2004.

[18] D. D. Djayaprawira, K. Tsunekawa, M. Nagai, H. Maehara, S. Yamagata, N. Watanabe, S. Yuasa, Y. Suzuki, and K. Ando, "230% room-temperature magnetoresistance in CoFeB/ MgO/ CoFeB magnetic tunnel junctions," *Applied Physics Letters*, vol. 86, no. 9, p. 092502, 2005.

[19] J. Hayakawa, S. Ikeda, F. Matsukura, H. Takahashi, and H. Ohno, "Dependence of giant tunnel magnetoresistance of sputtered CoFeB/MgO/CoFeB magnetic tunnel junctions on MgO barrier thickness and annealing temperature," *Japanese Journal of Applied Physics*, vol. 44, no. 4L, p. L587, 2005.

[20] S. Ikeda, J. Hayakawa, Y. Ashizawa, Y. Lee, K. Miura, H. Hasegawa, M. Tsunoda, F. Matsukura, and H. Ohno, "Tunnel magnetoresistance of 604% at 300 K by suppression of Ta diffusion in Co Fe B/ Mg O/ Co Fe B pseudo-spin-valves annealed at high temperature," *Applied Physics Letters*, vol. 93, no. 8, p. 082508, 2008.

[21] I. Prejbeanu, M. Kerekes, R. Sousa, H. Sibuet, O. Redon, B. Dieny, and J. Nozières, "Thermally assisted MRAM," *Journal of Physics: Condensed Matter*, vol. 19, no. 16, p. 165218, 2007.

[22] R. Bishnoi, F. Oboril, M. Ebrahimi, and M. B. Tahoori, "Self-timed read and write operations in STT-MRAM," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 24, no. 5, pp. 1783–1793, 2016.

[23] G. Prenat, K. Jabeur, P. Vanhauwaert, G. Di Pendina, F. Oboril, R. Bishnoi, M. Ebrahimi, N. Lamard, O. Boulle, K. Garello *et al.*, "Ultra-Fast and High-Reliability SOT-MRAM: From Cache Replacement to Normally-Off Computing," *IEEE Transactions on Multi-Scale Computing Systems*, vol. 2, no. 1, pp. 49–60, 2016.

[24] D. Suzuki, M. Natsui, A. Mochizuki, and T. Hanyu, "Cost-Efficient Self-Terminated Write Driver for Spin-Transfer-Torque RAM and Logic," *Transactions on Magnetics*, vol. 50, no. 11, pp. 1–4, 2014.

[25] S. Motaman, S. Ghosh, and N. Rathi, "Impact of process-variations in STTRAM and adaptive boosting for robustness," in *Design, Automation and Test in Europe Conference*, pp. 1431–1436, 2015.

[26] S. P. Mohanty, U. Choppali, and E. Kougianos, "Everything you wanted to know about smart cities: The Internet of things is the backbone," *IEEE Consumer Electronics Magazine*, vol. 5, no. 3, pp. 60–70, July 2016.

[27] D. Sabella, A. Vaillant, P. Kuure, U. Rauschenbach, and F. Giust, "Mobile-Edge Computing Architecture: The role of MEC in the Internet of Things," *IEEE Consumer Electronics Magazine*, vol. 5, no. 4, pp. 84–91, Oct 2016.

[28] J. Decuir, "The Story of the Internet of Things: Issues in utility, connectivity, and security." *IEEE Consumer Electronics Magazine*, vol. 4, no. 4, pp. 54–61, Oct 2015.

[29] D. A. Patterson and J. L. Hennessy, *Computer Organization and Design, Fourth Edition, Fourth Edition: The Hardware/Software Interface (The Morgan Kaufmann Series in Computer Architecture and Design)*, 4th ed. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2008.

[30] I. PRESENT, "Cramming more components onto integrated circuits," *Readings in computer architecture*, vol. 56, 2000.

[31] R. R. Schaller, "Moore's law: past, present and future," *IEEE Spectrum*, vol. 34, no. 6, pp. 52–59, Jun 1997.

[32] "Intel. [Online]." Available: http://www.intel.com/content/www/us/en/silicon-innovations/moores-law-technology.html, [accessed November 2016].

[33] S. Bob, "The Origin, Nature, and Implications of "MOORE'S LAW"," Available: http://research.microsoft.com/en-us/um/people/gray/Moore_Law.html, 1996, [accessed November 2016].

[34] W. A. Wulf and S. A. McKee, "Hitting the memory wall: implications of the obvious," *ACM SIGARCH computer architecture news*, vol. 23, no. 1, pp. 20–24, 1995.

[35] P. Machanick, "Approaches to addressing the memory wall," *School of IT and Electrical Engineering, University of Queensland*, 2002.

[36] M. Karunaratne and B. Oomann, "An optimal memory BISR implementation," *Journal of Computers*, vol. 8, no. 9, pp. 2167–2174, 2013.

[37] T. N. Theis and H.-S. P. Wong, "The End of Moore's Law: A New Beginning for Information Technology," *Computing in Science and Engineering*, vol. 19, no. 2, pp. 41–50, 2017.

[38] M. LaPedus, "What's Next For DRAM?" Available: http://semiengineering.com/whats-next-for-dram/, 2016, [accessed November 2016].

[39] K. Agarwal and S. Nassif, "The Impact of Random Device Variation on SRAM Cell Stability in Sub-90-nm CMOS Technologies," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 16, no. 1, pp. 86–97, Jan 2008.

[40] International Technology Roadmap for Semiconductors, Available: http://www.semiconductors.org/clientuploads/Research_Technology/ITRS/2015/0_2015%20ITRS%202.0%20Executive%20Report%20(1).pdf, 2015.

[41] S. Wolf, D. Awschalom, R. Buhrman, J. Daughton, S. Von Molnar, M. Roukes, A. Y. Chtchelkanova, and D. Treger, "Spintronics: a spin-based electronics vision for the future," *Science*, vol. 294, no. 5546, pp. 1488–1495, 2001.

[42] A. Fert *et al.*, "The origin, development and future of spintronics," *Phys. Usp*, vol. 178, no. 12, pp. 1336–1348, 2008.

[43] International Technology Roadmap for Semiconductors, Available: http://www.itrs.net, 2013.

[44] P. Chi, S. Li, Y. Cheng, Y. Lu, S. H. Kang, and Y. Xie, "Architecture design with STT-RAM: Opportunities and challenges," in *Asia and South Pacific Design Automation Conference*, pp. 109–114, 2016.

[45] Y. Zhang, X. Wang, Y. Li, A. K. Jones, and Y. Chen, "Asymmetry of MTJ switching and its implication to STT-RAM designs," in *Design, Automation and Test in Europe*, pp. 1313–1318, 2012.

[46] G. Sun, Y. Zhang, Y. Wang, and Y. Chen, "Improving energy efficiency of write-asymmetric memories by log style write," in *International Symposium on Low Power Electronics and Design*, pp. 173–178, 2012.

[47] X. Fong, S. H. Choday, and K. Roy, "Bit-cell level optimization for non-volatile memories using magnetic tunnel junctions and spin-transfer torque switching," *Transactions on Nanotechnology*, vol. 11, no. 1, pp. 172–181, 2012.

[48] J. Park, D. Ralph, and R. Buhrman, "Fast deterministic switching in orthogonal spin torque devices via the control of the relative spin polarizations," *Applied Physics Letters*, vol. 103, no. 25, p. 252406, 2013.

[49] K. Munira, W. H. Butler, and A. W. Ghosh, "A Quasi-Analytical Model for Energy-Delay-Reliability Tradeoff Studies During Write Operations in a Perpendicular STT-RAM Cell," *Transactions on Electron Devices*, vol. 59, no. 8, pp. 2221–2226, 2012.

[50] H. Kuriyama, T. Hirose, S. Murakami, T. Wada, K. Fujita, Y. Nishimura, and K. Anami, "An 8 ns 4 Mb serial access memory," *IEEE journal of solid-state circuits*, vol. 26, no. 4, pp. 502–506, 1991.

[51] C. S. Lytle and D. F. Faria, "Programmable logic array integrated circuit with general-purpose memory configurable as a random access or FIFO memory," Nov. 5 1996, US Patent 5,572,148.

[52] A. G. Fraser, "First-in, first-out (FIFO) memory configuration for queue storage," Mar. 26 1985, US Patent 4,507,760.

[53] K. Pagiamtzis and A. Sheikholeslami, "Content-addressable memory (CAM) circuits and architectures: A tutorial and survey," *IEEE Journal of Solid-State Circuits*, vol. 41, no. 3, pp. 712–727, 2006.

[54] K. Pagiamtzis and A. Sheikholeslami, "A low-power content-addressable memory (CAM) using pipelined hierarchical search scheme," *IEEE Journal of Solid-State Circuits*, vol. 39, no. 9, pp. 1512–1519, 2004.

[55] B. Keeth, R. J. Baker, , B. Johnson, and F. Lin, *DRAM Circuit Design: Fundamental and High-Speed Topics*, 2007.

[56] R. D. Adams, *High performance memory testing: design principles, fault modeling and self-test*, 2003.

*Bibliography*

[57] "ROM, EPROM, AND EEPROM TECHNOLOGY [Online]," Available: https://web.eecs.umich.edu/~prabal/teaching/eecs373-f10/readings/rom-eprom-eeprom-technology.pdf, [accessed Mar 2017].

[58] "microchip datasheet. [Online]." Available: https://www.jaapsch.net/psion/pdffiles/Eprom032k_datasheet_27C256.pdf, [accessed Mar 2017].

[59] "ST datasheet. [Online]." Available: http://www.st.com/content/ccc/resource/technical/document/datasheet/5c/df/52/a5/15/f2/48/bd/CD00259166.pdf/files/CD00259166.pdf/jcr:content/translations/en.CD00259166.pdf, [accessed Mar 2017].

[60] " [Online]." Available: http://www.eetimes.com/document.asp?doc_id=1272118, [accessed Mar 2017].

[61] A. Tal, "Two flash technologies compared: NOR vs NAND," *White Paper of M-SYstems*, 2002.

[62] F. Masuoka, M. Momodomi, Y. Iwata, and R. Shirota, "New ultra high density EPROM and flash EEPROM with NAND structure cell," in *Electron Devices Meeting*, pp. 552–555, 1987.

[63] S. Mittal and J. S. Vetter, "A Survey of Software Techniques for Using Non-Volatile Memories for Storage and Main Memory Systems," *IEEE Transactions on Parallel and Distributed Systems*, vol. 27, no. 5, pp. 1537–1550, 2016.

[64] B. prince, *Emerging Memories-Technologies and Trends*, 2002.

[65] R. Bez, E. Camerlenghi, A. Modelli, and A. Visconti, "Introduction to flash memory," *Proceedings of the IEEE*, vol. 91, no. 4, pp. 489–502, 2003.

[66] R. C. Sousa and I. L. Prejbeanu, "Non-volatile magnetic random access memories (MRAM)," *Comptes Rendus Physique*, vol. 6, no. 9, pp. 1013–1021, 2005.

[67] S. Tehrani, J. Slaughter, E. Chen, M. Durlam, J. Shi, and M. DeHerren, "Progress and outlook for MRAM technology," *IEEE Transactions on Magnetics*, vol. 35, no. 5, pp. 2814–2819, 1999.

[68] F. Oboril, F. Hameed, R. Bishnoi, A. Ahari, H. Naeimi, and M. Tahoori, "Normally-OFF STT-MRAM Cache with Zero-Byte Compression for Energy Efficient Last-Level Caches," in *Proceedings of International Symposium on Low Power Electronics and Design*, pp. 236–241, 2016.

[69] H. S. P. Wong, S. Raoux, S. Kim, J. Liang, J. P. Reifenberg, B. Rajendran, M. Asheghi, and K. E. Goodson, "Phase Change Memory," *Proceedings of the IEEE*, vol. 98, no. 12, pp. 2201–2227, 2010.

[70] S. Lai, "Current status of the phase change memory and its future," in *IEEE Electron Devices Meeting*, pp. 10–1, 2003.

[71] B. C. Lee, E. Ipek, O. Mutlu, and D. Burger, "Architecting phase change memory as a scalable dram alternative," in *ACM SIGARCH Computer Architecture News*, vol. 37, no. 3, pp. 2–13, 2009.

[72] L. Chua, "Resistance switching memories are memristors," in *Memristor Networks*. Springer, 2014, pp. 21–51.

[73] I. Baek, M. Lee, S. Seo, M. Lee, D. Seo, D.-S. Suh, J. Park, S. Park, H. Kim, I. Yoo *et al.*, "Highly scalable nonvolatile resistive memory using simple binary oxide driven by asymmetric unipolar voltage pulses," in *IEEE Electron Devices Meeting*, pp. 587–590, 2004.

[74] J. J. Yang, D. B. Strukov, and D. R. Stewart, "Memristive devices for computing," *Nature nanotechnology*, vol. 8, no. 1, pp. 13–24, 2013.

[75] J. F. Scott, P. De Araujo, and A. Carlos, "Ferroelectric memories," *Science(Washington, D. C.)*, vol. 246, no. 4936, pp. 1400–5, 1989.

[76] L. M. Grupp, J. D. Davis, and S. Swanson, "The bleak future of nand flash memory," in *Proceedings of the 10th USENIX conference on File and Storage Technologies*, pp. 2–2, 2012.

[77] L. Mearian, " Is NAND flash memory a dying technology ? [Online]." Available: IsNANDflashmemoryadyingtechnology, 2010, [accessed Mar 2017].

[78] A. Gebregiorgis, S. Kiamehr, F. Oboril, R. Bishnoi, and M. B. Tahoori, "A cross-layer analysis of Soft Error, aging and process variation in Near Threshold Computing," in *Design, Automation and Test in Europe (DATE)*, pp. 205–210, 2016.

[79] M. Ebrahimi, H. Asadi, R. Bishnoi, and M. B. Tahoori, "Layout-based modeling and mitigation of multiple event transients," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 35, no. 3, pp. 367–379, 2016.

[80] W. Thomson, "On the electro-dynamic qualities of metals:–effects of magnetization on the electric conductivity of nickel and of iron," *Proceedings of the Royal Society of London*, vol. 8, pp. 546–550, 1856.

[81] A. Fert, "Nobel Lecture: Origin, development, and future of spintronics," *Rev. Mod. Phys.*, vol. 80, pp. 1517–1530, Dec 2008.

[82] F. Mott, "Proc. R. Soc. London, Ser. A 153, 699," 1936.

[83] L. Berger, "Low-field magnetoresistance and domain drag in ferromagnets," *Journal of Applied Physics*, vol. 49, no. 3, pp. 2156–2161, 1978.

[84] P. Grunberg, "Magnetic field sensor with ferromagnetic thin layers having magnetically antiparallel polarized components," August 14 1990, US Patent 4,949,039.

[85] B. Dieny, "Giant magnetoresistance in spin-valve multilayers," *Journal of Magnetism and Magnetic Materials*, vol. 136, no. 3, pp. 335 – 359, 1994.

[86] G. Binasch, P. Grünberg, F. Saurenbach, and W. Zinn, "Enhanced magnetoresistance in layered magnetic structures with antiferromagnetic interlayer exchange," *Phys. Rev. B*, vol. 39, pp. 4828–4830, Mar 1989.

[87] M. Julliere, "Tunneling between ferromagnetic films," *Physics letters A*, vol. 54, no. 3, pp. 225–226, 1975.

[88] G. A. Prinz, "Magnetoelectronics," *Science*, vol. 282, no. 5394, pp. 1660–1663, 1998.

[89] E. Grochowski, "Emerging trends in data storage on magnetic hard disk drives," *Datatech (September 1998), ICG Publishing*, pp. 11–16, 1998.

[90] I. McFadyen, E. Fullerton, and M. Carey, "State-of-the-art magnetic hard disk drives," *Mrs Bulletin*, vol. 31, no. 05, pp. 379–383, 2006.

[91] C. Chappert, A. Fert, and F. N. Van Dau, "The emergence of spin electronics in data storage," *Nature materials*, vol. 6, no. 11, pp. 813–823, 2007.

[92] J. Daughton, "GMR applications," *Journal of Magnetism and Magnetic Materials*, vol. 192, no. 2, pp. 334–342, 1999.

[93] J. Daughton and Y. Chen, "GMR materials for low field applications," *IEEE Transactions on Magnetics*, vol. 29, no. 6, pp. 2705–2710, 1993.

[94] J.-G. J. Zhu and C. Park, "Magnetic tunnel junctions," *Materials Today*, vol. 9, no. 11, pp. 36–45, 2006.

[95] D. Monsma and S. Parkin, "Spin polarization of tunneling current from ferromagnet/Al2O3 interfaces using copper-doped aluminum superconducting films," *Applied Physics Letters*, vol. 77, no. 5, pp. 720–722, 2000.

[96] J. S. Moodera, L. R. Kinder, T. M. Wong, and R. Meservey, "Large magnetoresistance at room temperature in ferromagnetic thin film tunnel junctions," *Physical review letters*, vol. 74, no. 16, p. 3273, 1995.

[97] T. Miyazaki and N. Tezuka, "Giant magnetic tunneling effect in Fe/Al2O3/Fe junction," *Journal of Magnetism and Magnetic Materials*, vol. 139, no. 3, pp. L231–L234, 1995.

[98] W. Butler, X.-G. Zhang, T. Schulthess, and J. MacLaren, "Spin-dependent tunneling conductance of Fe| MgO| Fe sandwiches," *Physical Review B*, vol. 63, no. 5, p. 054416, 2001.

[99] B. Engel, J. Akerman, B. Butcher, R. Dave, M. DeHerrera, M. Durlam, G. Grynkewich, J. Janesky, S. Pietambaram, N. Rizzo *et al.*, "A 4-Mb toggle MRAM based on a novel bit and switching method," *IEEE Transactions on Magnetics*, vol. 41, no. 1, pp. 132–136, 2005.

*Bibliography*

[100] L. Savtchenko, B. Engel, N. Rizzo, M. Deherrera, and J. Janesky, "Method of writing to scalable magnetoresistance random access memory element," 2003, US Patent 6,545,906. [Online]. Available: https://www.google.ch/patents/US6545906

[101] M. El Baraji, V. Javerliac, W. Guo, G. Prenat, and B. Dieny, "Dynamic compact model of thermally assisted switching magnetic tunnel junctions," *Journal of Applied Physics*, vol. 106, no. 12, p. 123906, 2009.

[102] Y. Guillemenet, L. Torres, G. Sassatelli, N. Bruchon, and I. Hassoune, "A non-volatile run-time FPGA using thermally assisted switching MRAMS," in *International Conference on Field Programmable Logic and Applications*, pp. 421–426, 2008.

[103] I. Prejbeanu, S. Bandiera, J. Alvarez-Hérault, R. Sousa, B. Dieny, and J. Nozieres, "Thermally assisted MRAMs: ultimate scalability and logic functionalities," *Journal of Physics D: Applied Physics*, vol. 46, no. 7, p. 074002, 2013.

[104] J. E. Hirsch, "Spin Hall Effect," *Phys. Rev. Lett.*, vol. 83, pp. 1834–1837, Aug 1999.

[105] B. Miao, S. Huang, D. Qu, and C. Chien, "Inverse Spin Hall Effect in a Ferromagnetic Metal," *Physical review letters*, vol. 111, no. 6, p. 066602, 2013.

[106] "Wikipedia. [Online]." Available: https://en.wikipedia.org/wiki/Spin_Hall_effect, [accessed February 2017].

[107] T. Kawahara, K. Ito, R. Takemura, and H. Ohno, "Spin-transfer torque RAM technology: review and prospect," *Microelectronics Reliability*, vol. 52, no. 4, pp. 613–627, 2012.

[108] W. Guo, G. Prenat, V. Javerliac, M. E. Baraji, N. de Mestier, C. Baraduc, and B. Dieny, "SPICE modelling of magnetic tunnel junctions written by spin-transfer torque," *Journal of Applied Physics*, vol. 43, no. 21, p. 215001, May 2010.

[109] D. Apalkov, A. Khvalkovskiy, S. Watts, V. Nikitin, X. Tang, D. Lottis, K. Moon, X. Luo, E. Chen, A. Ong, A. Driskill-Smith, and M. Krounbi, "Spin-transfer torque magnetic random access memory (STT-MRAM)," *Journal of Emerging Technologies in Computer Systems*, vol. 9, no. 2, pp. 13:1–13:35, 2013.

[110] N. Sayed, F. Oboril, R. Bishnoi, and M. B. Tahoori, "Leveraging Systematic Unidirectional Error-Detecting Codes for fast STT-MRAM cache," in *VLSI Test Symposium (VTS)*, pp. 1–6, 2017.

[111] N. Sayed, M. Ebrahimi, R. Bishnoi, and M. B. Tahoori, "Opportunistic write for fast and reliable STT-MRAM," in *Design, Automation Test in Europe Conference Exhibition (DATE)*, pp. 554–559, 2017.

[112] C. Yoshida and T. Sugii, "Reliability study of magnetic tunnel junction with naturally oxidized MgO barrier," in *International Reliability Physics Symposium*, pp. 2A–3, 2012.

[113] C. Yoshida, M. Kurasawa, Y. M. Lee, K. Tsunoda, M. Aoki, and Y. Sugiyama, "A study of dielectric breakdown mechanism in CoFeB/MgO/CoFeB magnetic tunnel junction," in *International Reliability Physics Symposium*, pp. 139–142, 2009.

[114] G. Panagopoulos, C. Augustine, and K. Roy, "Modeling of dielectric breakdown-induced time-dependent STT-MRAM performance degradation," in *Device Research Conference*, pp. 125–126, 2011.

[115] S. M. Nair, R. Bishnoi, M. S. Golanbari, F. Oboril, and M. B. Tahoori, "VAET-STT: A variation aware estimator tool for STT-MRAM based memories," in *Design, Automation Test in Europe Conference Exhibition (DATE)*, pp. 1456–1461, 2017.

[116] A. Ahari, M. Ebrahimi, F. Oboril, and M. Tahoori, "Improving reliability, performance, and energy efficiency of STT-MRAM with dynamic write latency," in *IEEE International Conference on Computer Design (ICCD)*, pp. 109–116, 2015.

[117] R. Robertazzi, J. Nowak, and J. Sun, "Analytical MRAM test," in *ITC*, pp. 1–10, 2014.

[118] K. Sugiura, S. Takahashi, M. Amano, T. Kajiyama, M. Iwayama, Y. Asao, N. Shimomura, T. Kishi, S. Ikegawa, H. Yoda *et al.*, "Ion beam etching technology for high-density spin transfer torque magnetic random access memory," *JJAP*, vol. 48, no. 8S1, p. 08HD02, 2009.

[119] H. Naemi, C. Augustine, A. Raychowdhury, S. Lu, J. Tschanz, "STTRAM Scaling And Retention Failure," *Intel Technology Journal*, vol. 17, no. 1, pp. 54–75, 2013.

[120] M. Kuepferling, S. Zullino, A. Sola, B. Van de Wiele, G. Durin, M. Pasquale, K. Rott, G. Reiss, and G. Bertotti, "Vortex dynamics in Co-Fe-B magnetic tunnel junctions in presence of defects," *JAP*, vol. 117, no. 17, p. 17E107, 2015.

[121] J. Sun, M. Gaidis, G. Hu, E. O'Sullivan, S. Brown, J. Nowak, P. Trouilloud, and D. Worledge, "High-bias backhopping in nanosecond time-domain spin-torque switches of MgO-based magnetic tunnel junctions," *JAP*, vol. 105, no. 7, p. 07D109, 2009.

[122] T. Min, J. Sun, R. Beach, D. Tang, and P. Wang, "Back-hopping after spin torque transfer induced magnetization switching in magnetic tunneling junction cells," *JAP*, vol. 105, no. 7, p. 07D126, 2009.

[123] K. Ishida, T. Yasufuku, S. Miyamoto, H. Nakai, M. Takamiya, T. Sakurai, and K. Takeuchi, "A 1.8V 30nJ adaptive program-voltage (20V) generator for 3D-integrated NAND flash SSD," in *IEEE International Solid-State Circuits Conference - Digest of Technical Papers*, pp. 238–239,239a, 2009.

[124] F. Oboril, R. Bishnoi, M. Ebrahimi, and M. B. Tahoori, "Evaluation of hybrid memory technologies using SOT-MRAM for on-chip cache hierarchy," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 34, no. 3, pp. 367–380, 2015.

[125] R. Bishnoi, M. Ebrahimi, F. Oboril, and M. B. Tahoori, "Architectural aspects in design and analysis of SOT-based memories," in *Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 700–707, 2014.

[126] F. Oboril, R. Bishnoi, M. Ebrahimi, M. Tahoori, G. Di Pendina, K. Jabeur, and G. Prenat, "Spin Orbit Torque memory for non-volatile microprocessor caches," in *International Workshop on Emerging Memory Solutions, DATE*, 2016.

[127] S. Mittal, R. Bishnoi, F. Oboril, H. Wang, M. Tahoori, A. Jog, and J. Vetter, "Architecting SOT-RAM Based GPU Register File," in *Computer Society Annual Symposium on VLSI (ISVLSI)*, 2017.

[128] M. B. Tahoori and S.M. Nair and R. Bishnoi and S. Senni and J. Mohdad and F. Mailly and L. Torres and P. Benoit and P. Nouet and R. Ma and M. Kreißig and F. Ellinger and K. Jabeur and P. Vanhauwaert and G. Di Pendina and G. Prenat, "GREAT: heteroGeneous integRated magnetic tEchnology using multifunctional standardized sTack," in *Computer Society Annual Symposium on VLSI (ISVLSI)*, 2017.

[129] X. Dong, C. Xu, Y. Xie, and N. P. Jouppi, "Nvsim: A circuit-level performance, energy, and area model for emerging nonvolatile memory," *Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 31, no. 7, pp. 994–1007, 2012.

[130] N. Muralimanohar, R. Balasubramonian, and N. P. Jouppi, "CACTI 6.0: A tool to model large caches," *HP Laboratories*, pp. 22–31, 2009.

[131] A. Raychowdhury, "Pulsed READ in spin transfer torque (STT) memory bitcell for lower READ disturb," in *International Symposium on Nanoscale Architectures (NANOARCH)*, pp. 34–35, 2013.

[132] P. Zhou, B. Zhao, J. Yang, and Y. Zhang, "Energy reduction for STT-RAM using early write termination," in *International Conference on Computer-Aided Design-Digest of Technical Papers*, pp. 264–268, 2009.

[133] R. Bishnoi, F. Oboril, M. Ebrahimi, and M. Tahoori, "Avoiding Unnecessary Write Operations in STT-MRAM for Low Power Implementation," in *International Symposium on Quality Electronic Design*, pp. 548–553, 2014.

[134] X. Bi, Z. Sun, H. Li, and W. Wu, "Probabilistic design methodology to improve run-time stability and performance of STT-RAM caches," in *International Conference on Computer-Aided Design-Digest of Technical Papers*, pp. 88–94, 2012.

[135] R. Bishnoi, M. Ebrahimi, F. Oboril, and M. Tahoori, "Asynchronous asymmetrical write termination (AAWT) for a low power STT-MRAM," in *Design, Automation and Test in Europe*, pp. 1–6, Mar. 2014.

[136] T. Zheng, J. Park, M. Orshansky, and M. Erez, "Variable-energy write STT-RAM architecture with bit-wise write-completion monitoring," in *International Symposium on Low Power Electronics and Design*, pp. 229–234, Sept 2013.

# Bibliography

[137] V. W.-Y. Sit, C.-S. Choy, and C.-F. Chan, "A four-phase handshaking asynchronous static RAM design for self-timed systems," *Journal of Solid-State Circuits*, vol. 34, no. 1, pp. 90–96, 1999.

[138] V. N. Ekanayake and R. Manohar, "Asynchronous DRAM design and synthesis," in *International Symposium on Asynchronous Circuits and Systems*, pp. 174–183, 2003.

[139] B. S. Amrutur and M. A. Horowitz, "A replica technique for wordline and sense control in low-power SRAM's," *Journal of Solid-State Circuits*, vol. 33, no. 8, pp. 1208–1219, 1998.

[140] M.-F. Chang, S.-M. Yang, and K.-T. Chen, "Wide VDD embedded asynchronous SRAM with dual-mode self-timed technique for dynamic voltage systems," *Transactions on Circuits and Systems Part I: Regular Papers*, vol. 56, no. 8, pp. 1657–1667, 2009.

[141] W. Zhao, C. Chappert, V. Javerliac, and J.-P. Noziere, "High speed, high stability and low power sensing amplifier for MTJ/CMOS hybrid logic circuits," *Transactions on Magnetics*, vol. 45, no. 10, pp. 3784–3787, 2009.

[142] M. Margala, "Low-power SRAM circuit design," in *International Workshop on Memory Technology, Design and Testing*, pp. 115–122, 1999.

[143] R. Fengbo, H. Park, R. Dorrance, Y. Toriyama, C.-K. Yang, and D. Markovic, "A body-voltage-sensing-based short pulse reading circuit for spin-torque transfer RAMs (STT-RAMs)," in *International Symposium on Quality Electronic Design*, pp. 275–282, 2012.

[144] Y. Umeki, K. Yanagida, S. Yoshimoto, S. Izumi, M. Yoshimoto, H. Kawaguchi, K. Tsunoda, and T. Sugii, "A negative-resistance sense amplifier for low-voltage operating STT-MRAM," in *Asia and South Pacific Design Automation Conference*, pp. 8–9, 2015.

[145] M. R. Guthaus, J. S. Ringenberg, D. Ernst, T. M. Austin, T. Mudge, and R. B. Brown, "MiBench: A free, commercially representative embedded benchmark suite," in *International Workshop on Workload Characterization*, pp. 3–14, 2001.

[146] R. Bishnoi, M. Ebrahimi, F. Oboril, and M. B. Tahoori, "Improving Write Performance for STT-MRAM," *IEEE Transactions on Magnetics*, vol. 52, no. 8, pp. 1–11, 2016.

[147] C. W. Smullen, V. Mohan, A. Nigam, S. Gurumurthi, and M. R. Stan, "Relaxing non-volatility for fast and energy-efficient STT-RAM caches," in *High Performance Computer Architecture*, pp. 50–61, 2011.

[148] A. Nigam, C. Smullen, V. Mohan, E. Chen, S. Gurumurthi, and M. R. Stan, "Delivering on the promise of universal memory for spin-transfer torque RAM (STT-RAM)," in *Proceedings of the International Symposium on Low Power Electronics and Design*, pp. 121–126, 2011.

[149] B. Del Bel, J. Kim, C. H. Kim, and S. S. Sapatnekar, "Improving STT-MRAM density through multibit error correction," in *Design, Automation and Test in Europe*, pp. 1–6, 2014.

[150] N. Sayed, F. Oboril, A. Shirvanian, R. Bishnoi, and M. B. Tahoori, "Exploiting STT-MRAM for Approximate Computing," in *European Test Symposium (ETS)*, 2017.

[151] D. Lee, S. K. Gupta, and K. Roy, "High-performance low-energy STT MRAM based on balanced write scheme," in *International Symposium on Low power Electronics and Design*, pp. 9–14, 2012.

[152] R. Patel, X. Guo, Q. Guo, E. Ipek, and E. Friedman, "Reducing Switching Latency and Energy in STT-MRAM Caches With Field-Assisted Writing," *Transactions on Very Large Scale Integration (VLSI) Systems*, vol. PP, no. 99, pp. 1–1, 2015.

[153] X. Fong, Y. Kim, S. Choday, and K. Roy, "Failure Mitigation Techniques for 1T-1MTJ Spin-Transfer Torque MRAM Bit-cells," *Transactions on Very Large Scale Integration Systems*, vol. 22, no. 2, pp. 384–395, Feb 2014.

[154] Y. Kim, S. K. Gupta, S. P. Park, G. Panagopoulos, and K. Roy, "Write-optimized reliable design of STT MRAM," in *International Symposium on Low power Electronics and Design*, pp. 3–8, 2012.

[155] J. Li, P. Ndai, A. Goel, S. Salahuddin, and K. Roy, "Design paradigm for robust spin-torque transfer magnetic RAM (STT MRAM) from circuit/architecture perspective," *Transactions on Very Large Scale Integration Systems*, vol. 18, no. 12, pp. 1710–1723, 2010.

[156] R. Bishnoi, M. Ebrahimi, F. Oboril, and M. Tahoori, "Read disturb fault detection in STT-MRAM," in *International Test Conference*, pp. 1–7, Oct 2014.

[157] X. Wang, Z. Wang, X. Hao, Y. Zhou, J. Zhang, H. Gan, D. H. Jung, K. Satoh, B. Yen, R. Malmhall *et al.*, "Different dielectric breakdown mechanisms for RF-MgO and naturally oxidized MgO," *Applied Physics Express*, vol. 7, no. 8, p. 083002, 2014.

[158] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti *et al.*, "The gem5 simulator," *ACM SIGARCH Computer Architecture News*, vol. 39, no. 2, pp. 1–7, 2011.

[159] R. Bishnoi, F. Oboril, and M. B. Tahoori, "Low-power multi-port memory architecture based on Spin Orbit Torque magnetic devices," in *Great Lakes Symposium on VLSI*, pp. 409–414, 2016.

[160] Y. Seo, X. Fong, K.-W. Kwon, and K. Roy, "Spin-Hall Magnetic Random-Access Memory With Dual Read/Write Ports for On-Chip Caches," *Magnetics Letters*, vol. 6, pp. 1–4, 2015.

[161] X. Bi, M. A. Weldon, and H. Li, "STT-RAM designs supporting dual-port accesses," in *DATE*, pp. 853–858, 2013.

[162] X. Liu, M. Mao, X. Bi, H. Li, and Y. Chen, "An efficient STT-RAM-based register file in GPU architectures," *ASP-DAC*, pp. 490–495, 2015.

[163] K. Jabeur, G. Di Pendina, G. Prenat, L. Buda-Prejbeanu, and B. Dieny, "Compact Modeling of a Magnetic Tunnel Junction Based on Spin Orbit Torque," *Magnetics, IEEE Transactions on*, vol. 50, no. 7, pp. 1–8, 2014.

[164] K. Jabeur, G. Di Pendina, and G. Prenat, "Ultra-energy-efficient CMOS/magnetic nonvolatile flip-flop based on spin-orbit torque device," *EL*, vol. 50, no. 8, pp. 585–587, 2014.

[165] K. Jabeur, G. Di Pendina, F. Bernard-Granger, and G. Prenat, "Spin orbit torque non-volatile flip-flop for high speed and low energy applications," *Electron Device Letters*, vol. 35, no. 3, pp. 408–410, 2014.

[166] R. Bishnoi, F. Oboril, and M. B. Tahoori, "Non-Volatile Non-Shadow flip-flop using Spin Orbit Torque for efficient normally-off computing," in *Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 769–774, 2016.

[167] Y. Liu, Z. Chen, S. Zhang, and J. Zhou, "State retention power gated cell," Aug. 26 2014, US Patent 8,816,741.

[168] M. Padhye and D. Gross, "Freescale: Wireless Low-Power Design and Verification with CPF," https://www.si2.org/?page=1061.

[169] S. Yamamoto, Y. Shuto, and S. Sugahara, "Nonvolatile flip-flop using pseudo-spin-transistor architecture and its power-gating applications," in *ISCDG*, pp. 17–20, 2012.

[170] D. Chabi, W. Zhao, E. Deng, Y. Zhang, N. Ben Romdhane, J.-O. Klein, and C. Chappert, "Ultra Low Power Magnetic Flip-Flop Based on Checkpointing/Power Gating and Self-Enable Mechanisms," *Transactions on Circuits and Systems*, vol. 61, no. 6, pp. 1755–1765, 2014.

[171] K.-W. Kwon, S. H. Choday, Y. Kim, X. Fong, S. P. Park, and K. Roy, "SHE-NVFF: spin Hall effect-based nonvolatile flip-flop for power gating architecture," *Electron Device Letters*, vol. 35, no. 4, pp. 488–490, 2014.

[172] Garello, C. O. Avci, I. M. Miron, M. Baumgartner, A. Ghosh, S. Auffret, O. Boulle, G. Gaudin, and P. Gambardella, "Ultrafast magnetization switching by spin-orbit torques," *Applied Physics Letters*, vol. 105, no. 21, 2014.

[173] Hanson, S. and Mingoo Seok and Yu-Shiang Lin and Zhiyoong Foo and Daeyeon Kim and Yoonmyung Lee and Liu, N. and Sylvester, D. and Blaauw, D., "A Low-Voltage Processor for Sensing Applications With Picowatt Standby Mode," *Solid-State Circuits*, vol. 44, no. 4, pp. 1145–1155, 2009.

[174] N. Choudhary, S. Wadhavkar, T. Shah, H. Mayukh, J. Gandhi, B. Dwiel, S. Navada, H. Najaf-abadi, and E. Rotenberg, "FabScalar: Automating Superscalar Core Design," *Micro*, vol. 32, no. 3, pp. 48–59, 2012.

[175] R. Bishnoi, F. Oboril, and M. B. Tahoori, "Fault tolerant non-volatile spintronic flip-flop," in *Design, Automation and Test in Europe Conference (DATE)*, pp. 261–264, 2016.

[176] R. Bishnoi, F. Oboril, and M. B. Tahoori, "Design of Defect and Fault-Tolerant Nonvolatile Spintronic Flip-Flops," *IEEE Transactions on Very Large Scale Integration (TVLSI)*, vol. 25, no. 4, pp. 1421–1432, 2017.

[177] W. Zhao, E. Belhaire, and C. Chappert, "Spin-MTJ based non-volatile flip-flop," in *NANO*, pp. 399–402, 2007.

[178] I. Kim, Y. Zorian, G. Komoriya, H. Pham, F. P. Higgins, and J. L. Lewandowski, "Built in self repair for embedded high density SRAM," in *ITC*, pp. 1112–1119, Oct 1998.

[179] J.-J. Shau, "Methods to make DRAM fully compatible with SRAM using error correction code (ECC) mechanism," Apr. 10 2001, US Patent 6,216,246.

[180] W. Liu, J. Rho, and W. Sung, "Low-power high-throughput BCH error correction VLSI design for multi-level cell NAND flash memories," in *IEEE Workshop on Signal Processing Systems Design and Implementation*, pp. 303–308, 2006.

[181] A. M. Yamauchi, "Error detection and correction code for data and check code fields," Mar. 21 2000, uS Patent 6,041,430.

[182] S. Khanna, S. C. Bartling, M. Clinton, S. Summerfelt, J. A. Rodriguez, and H. P. McAdams, "An FRAM-Based Nonvolatile Logic MCU SoC Exhibiting 100Digital State Retention at VDD= 0 V Achieving Zero Leakage With < 400-ns Wakeup Time for ULP Applications," *Journal of Solid-State Circuits*, vol. 49, no. 1, pp. 95–106, 2014.

[183] W. Kang, W. Zhao, Z. Wang, Y. Zhang, J. O. Klein, C. Chappert, Y. Zhang, and D. Ravelosona, "DFSTT-MRAM: Dual Functional STT-MRAM Cell Structure for Reliability Enhancement and 3-D MLC Functionality," *Transactions on Magnetics*, vol. 50, no. 6, pp. 1–7, 2014.

[184] Y. Tsuji, R. Nebashi, N. Sakimura, A. Morioka, H. Honjo, K. Tokutome, S. Miura, T. Suzuki, S. Fukami, K. Kinoshita, T. Hanyu, T. Endoh, N. Kasai, H. Ohno, and T. Sugibayashi, "Spintronics primitive gate with high error correction efficiency 6(Perror)2 for logic-in memory architecture," in *Symposium on VLSI Technology*, pp. 63–64, 2012.

[185] W. kang, W. Zhao, E. Deng, J.-O. Klein, Y. Cheng, D. Ravelosona, Y. Zhang, and C. Chappert, "A radiation hardened hybrid spintronic/CMOS nonvolatile unit using magnetic tunnel junctions," *Journal of Physics D: Applied Physics*, vol. 47, no. 40, p. 405003, 2014.

[186] A. Mejdoubi, G. Prenat, and B. Dieny, "A compact model of precessional spin-transfer switching for MTJ with a perpendicular polarizer," in *MIEL*, pp. 225–228, 2012.

[187] G. Prenat, B. Dieny, J. Nozieres, G. DiPendina, and K. Torki, "Hybrid CMOS/Magnetic Process Design Kit and application to the design of high-performances non-volatile logic circuits," in *ICCAD*, pp. 240–245, 2011.