

MAGNETIC COUPLED SPIN-TORQUE DEVICE:
SPIN-BASED NON-VOLATILE LOGIC DEVICE AND
APPLICATIONS

A DISSERTATION
SUBMITTED TO THE DEPARTMENT OF
ELECTRICAL ENGINEERING
AND THE COMMITTEE ON GRADUATE STUDIES
OF STANFORD UNIVERSITY
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

Larkhoon Leem

August 2010

© 2010 by Lark-Hoon Leem. All Rights Reserved.

Re-distributed by Stanford University under license with the author.



This work is licensed under a Creative Commons Attribution-Noncommercial 3.0 United States License.

<http://creativecommons.org/licenses/by-nc/3.0/us/>

This dissertation is online at: <http://purl.stanford.edu/xt251xv8132>

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

James Harris, Primary Adviser

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

Peter Peumans

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

Robert M White

Approved for the Stanford University Committee on Graduate Studies.

Patricia J. Gumport, Vice Provost Graduate Education

This signature page was generated electronically upon submission of this dissertation in electronic format. An original signed hard copy of the signature page is on file in University Archives.

Abstract

Power consumption has become the key constraint in electronics design, since the MOSFET threshold voltage (V_{th}) and hence the supply voltage (V_{dd}) can no longer be scaled. This trend calls for new device concepts such as Spintronic devices that are fundamentally different from CMOS. However, the MOSFET-type Spintronics transistor has not been demonstrated yet due to the technical difficulties in injecting, transporting and detecting spin information. In this work, I present an alternative Spintronics logic device, *Magnetic Coupled Spin Torque Device (MCSTD)*, which is free from spin-injection, transport and detection problems. It leverages spin-torque transfer effect and magnetic dipole coupling between spin-torque devices to modulate its magnetization reversal energy barrier. Its device switching speed, signal inversion and signal level restoration capabilities will be discussed. For device-to-device level spin communication, MCSTD uses a novel interconnection technique that efficiently converts spin (or magneto-resistance) information to current amplitude difference information, which is then converted back to spin information at the subsequent gates. In micro-magnetic simulations, MCSTD-based NAND, NOR, XOR gates and a three-stage ring oscillator have been demonstrated to estimate realistic device speed and power consumption. The fabrication of 20nm gap MCSTDs has been successfully completed and demonstrated the input dependent switching voltage modulation, i.e., switching voltage of MCSTD gate depends on the magnetic orientations of the input spin-torque devices. The amount of voltage shifts ranged between 40~300mV, which is well above thermal fluctuations. Non-volatility in logic device such as MCSTD opens up very unique potential applications in future

power management techniques and smart sensor technologies. For example, MCSTDs can replace SRAMs and pass gate transistors in reconfigurable logics such as Field Programmable Gate Array (FPGA). Instant-on/off nature of MCSTD enables low overhead system-level power gating scheme for embedded devices. Also, MCSTD can be used as a magnetic sensor with in-situ logic operations for error-resilient DNA microarray sensors.

This work also explored other disruptive low-power device and system solutions such as Graphene nano-ribbon /Carbon nanotube based heterojunction Tunneling FET (Chap.8) and Error Resilient System Architecture for Probabilistic Applications (Chap.9).

Acknowledgments

As I finish my Ph.D., I realize how lucky I was at the turning points of my journey during my Ph.D. I thank all those people who stood by me and helped me make it through those difficult times. First of all, I would like to thank Prof. James S. Harris, my principle thesis advisor. Before I joined Harris group, I was coming from considerably distant field of research. Prof. Harris, *the Coach*, trusted me somehow and helped me have a fresh new start in his group. I am sure this was not an easy decision because, I did not have any proven research record in his area, which many other faculty members request. Ever since I worked with him, he always supported in full and let me be in the driver's seat regarding my Ph.D research. I am confident his trust on me made me to work even harder to prove that he was right. Now I realize that's his style of training students.

I thank Dr. Stuart S. P. Parkin at IBM Almaden research center for allowing me to collaborate with his world-leading research staff members and use his state-of-the-art facilities. Without the expertise on Spintronics and spin-torque device fabrication in his group, we would not have the working prototype MCSTD devices. Among IBM research staff members, I'd like to thank Dr. Xin Jiang, Dr. See-hun Yang, Dr. Charles Rettner and Brian Hughes for wonderful support and advice. They are simply the "best" people in spin-torque transfer technology.

I would like to thank my Ph.D Orals committee members. Prof. Peter Peumans is one of the "hardest-to-catch" professors. But, I was able to run my ideas on new inventions (ones that Prof. Harris didn't show interest) to him and he always gave me valuable feedbacks. I also thank Prof. Boris Murmann for being my Ph.D Orals chair. His feedbacks from circuit designer's perspective became very important part of findings in this work. I appreciate Prof. Robert M. White for sharing his insights on the topics of magnetism. He generously spared his time to have long discussions.

I thank all my friends in Harris group and Mitra group and Korean Students Associations.

Finally, I thank my wife, Kijung and my son, Ethan Dongha for being my endless source of joy. Since Kijung and I are both graduate students, raising a child is not an easy task. But the happiness of being together literally outweighs any troubles and it is as rewarding as it can be. My parents, Mr. Lae-gue Leem and Mrs. Duk-yee Kim have not retired yet, and I must be the primary reason. I thank them and love them.

I thank everybody with all my heart and wish happiness and well-being.

Table of Contents

Abstract	iii
Acknowledgments	v
List of Tables	x
List of Figures	xi
Chapter 1. Introduction	1
Chapter 2. Introduction to Spintronics and Non-volatile Applications	7
2.1 Motivation	8
2.2 Spin-orbit interaction	11
2.3 Spin dependent electronic structures	13
2.4 Magneto-resistance (MR)	15
2.5 Spin-torque transfer effect and magnetic reversal.....	18
2.6 Electric field based switching	25
2.7 Spintronic logic device	27
Chapter 3. Overview of Magnetic Coupled Spin-Torque Device (MCSTD)	31
3.1 Introduction	32
3.2 MCSTD Structure and Operation	32
3.3 Functional Logic Design	37
3.3.1 NAND and NOR gates	38
3.3.2 XOR and XNOR gates	40
3.4 Design considerations	42
3.4.1 Voltage shift	42
3.4.2 Output MTJ driver current and device reset	44

3.4.3 Oersted field induced noise	45
3.4.4 Positions and angles of the input spin-torque devices	46
3.4.5 Materials for the input device	49
3.5 MCSTD Logic Design	51
3.5.1 Gain	51
3.5.2 Nonlinearity	52
3.5.3 Cascadability	53
3.5.4 Fanout	54
3.5.5 Static power consumption	54
Chapter 4. Spin interconnection	57
4.1 Introduction	58
4.2 Complementary MCSTD based spin-interconnection	58
4.3 Magnetic Ring Oscillator	61
4.4 Magnetic domain-wall based interconnections	63
4.4.1 Overview	63
4.4.2 Reliable signal interconnection scheme	64
Chapter 5. Applications of Magnetic Coupled Spin-Torque Devices	68
5.1 Introduction	69
5.2 MCSTD based 2-bit full adder	69
5.2.1 Full adder design	69
5.2.2 Energy estimate	71
5.3 Logic embedded bio/image sensor	74
5.4 MCSTD for future reconfigurable logic	76
5.4.1 MCSTD in crossbar array architecture	76
5.4.2 Speedup of MCSTD over conventional MRAM technology	78
5.4.3 Incorporation of MCSTD into CMOS	82
5.4.4 MCSTD based Look-up Table (LUT)	84
5.4.5 MCSTD based Routing Fabric	86
5.4.6 MCSTD based Flip-flop	87

5.5 Simulation result of MCSTD for future FPGA architecture	88
5.5.1 MCSTD vs. CMOS	88
5.5.2 MCSTD based FPGA Look-up Table (LUT)	90
5.5.3 MCSTD based Routing Fabric	92
5.5.4 MCSTD based Flip-flop	93
Chapter 6. Fabrication process of Magnetic Coupled Spin-Torque Devices	97
6.1 Fabrication challenges	98
6.2 MTJ film stack	99
6.3 Bottom contact	99
6.4 E-beam lithography	101
6.5 Chrome hard mask layer	102
6.6 Ion-milling the gaps between MTJs	103
6.7 Effect of shape irregularities	105
6.8 Electrical contacts	105
Chapter 7. Experimental measurements of Magnetically Coupled Spin-Torque Devices (MCSTD)	108
7.1 Introduction	109
7.2 Samples and Measurement setup	109
7.3 Input-dependent switching voltage shift	111
7.4 Input dependent switching voltage shift vs. H_{noise} and H_C	113
7.5 Energy barrier height measurements & switching speed	118
7.6 Switching probability and logic operations	120
7.7 Energy consumption measurements	126
Chapter 8. Multi-scale Simulations of Partially Unzipped CNT Hetero-junction Tunneling Field Effect Transistor	130
8.1 Motivation	131
8.2 Device simulation process	132
8.3 Device operating principle and simulation results	133

8.4 Conclusion	138
Chapter 9. ERSA: Error Resilient System Architecture for Probabilistic Applications	140
9.1 Introduction	141
9.2 Error resilience of probabilistic applications	143
9.3 ERSA Overview	144
9.3.1. Super Reliable Core (SRC)	145
9.3.2. Relaxed-Reliability Cores (RRCs)	146
9.3.3. Algorithmic Convergence Test	147
9.3.3.1. Convergence Damping	147
9.3.3.2. Convergence Filtering	149
9.4. ERSA Experiments	149
9.4.1. ERSA Experimental Results: Logic Errors in RRCs	149
9.4.1.1. ERSA Computation Accuracy	151
9.4.1.2. ERSA Execution time	152
9.4.2. RRC L1 Data Cache Errors	153
Chapter 10. Conclusions and future works	158

List of Tables

Table 5.1:	Comparison of advanced memories	85
Table 5.2:	Device count comparison between CMOS and MCSTD	88
Table 5.3:	Gate area comparison between CMOS and MCSTD	88
Table 5.4:	Comparison of CMOS LUT and nano-magnet/CMOS hybrid circuit LUT energy consumption	90
Table 7.1:	Truth table for NAND and NOR logic	122
Table 8.1:	Electronic parameters calculated using DFT and EHT theories	133
Table 8.2:	Fabrication benefits of GNR/CNT Heterojunction TFET vs. GNR, CNT TFET	138

List of Figures

Figure 1.1:	Specint2006 (www.spec.org/cpu2006/results) benchmark results of modern CMOS microprocessors	2
Figure 1.2:	Power related metrics for MOSFET scaling	3
Figure 2.1:	Exchange coupling of electron spins (a) positive coupling (b) negative coupling.....	8
Figure 2.2:	Typical Multi-Aperture Device	9
Figure 2.3:	Magnetic Tunnel Junction	10
Figure 2.4:	Spin-Orbit Interaction induced magnetic fields	11
Figure 2.5:	Spin-dependent scattering and Spin-Hall effect	12
Figure 2.6:	Density of States (DOSs) comparison of magnetic (Co) and non-magnetic (Cu) materials	13
Figure 2.7:	Density of States (DOSs) comparison of magnetic tunnel junctions.....	14
Figure 2.8:	Evolution of magnetoresistive film structures	16
Figure 2.9:	Areal-density growth curve of Hard Disk Drive (HDD) recording products	16
Figure 2.10:	Current-in-plane (CIP) vs. Current-perpendicular-to-plane (CPP) readhead architecture	17
Figure 2.11:	Energy barrier of nano-pillar shaped spin-torque device.....	19
Figure 2.12:	Micromagnetics modeling of magnetic reversal.....	22
Figure 2.13:	Operation modes of spin-torque devices	23
Figure 2.14:	Spin-torque transfer switching of magnetic tunnel junction.....	24
Figure 2.15:	Spin current driven magnetization switching phase diagram.....	24
Figure 2.16:	Electric field based switching device structure.....	25
Figure 2.17:	Spin-based current gating of Spin-FET	27

Figure 3.1:	Energy barrier of nano-pillar shaped spin-torque device.....	33
Figure 3.2:	Schematic and SEM photo of our Magnetic Coupled Spin-torque Device (MCSTD).....	34
Figure 3.3:	Net magnetic fringing fields	35
Figure 3.4:	Example design of MCSTD and its energy barriers at different input MTJ magnetization configurations	36
Figure 3.5:	Simplified magnetic energy model of MCSTD device	36
Figure 3.6:	Input signal dependent energy barrier height.....	37
Figure 3.7:	MCSTD NAND and NOR gate concept	39
Figure 3.8:	MCSTD NAND gate energy barrier height versus input device magnetization directions	39
Figure 3.9:	Red : spin up, Blue : spin down, $J=3.63 \times 10^7$ A/cm ² , switching happens at (1,1)=(spin_up, spin_up) input only	40
Figure 3.10:	MCSTD XOR gate energy barrier height vs. the input MTJ magnetizations.....	41
Figure 3.11:	The shift of MR vs. H-field loop due to net fringing fields from the input MTJs	42
Figure 3.12:	Energy profiles of uni-stable state MCSTD and clocking schemes.....	44
Figure 3.13:	Influence of the Oersted field on spin-torque switching process.....	45
Figure 3.14:	Dipole coupling between the input and output MTJs.....	47
Figure 3.15:	Dipole coupling between two magnetic moments	47
Figure 3.16:	Dipole coupling among three magnetic moments	47
Figure 3.17:	MCSTD gate with the input MTJ at different locations and angles.....	48
Figure 3.18:	Dipole coupling strength vs. input MTJ angles	48
Figure 3.19:	Proposed voltage operation range of MCSTD gates	50
Figure 3.20:	MCSTD NAND gate energy barrier heights	50
Figure 3.21:	Horizontal components of fringing magnetic fields	51
Figure 3.22:	Circuit configuration for MCSTDs	53
Figure 4.1:	Complementary MCSTD and spin interconnection	59
Figure 4.2:	De Morgan's law	60

Figure 4.3:	Actual implementation of complementary MCSTD based spin interconnection	60
Figure 4.4:	A three stage MCSTD ring oscillator	61
Figure 4.5:	MCSTD ring oscillator frequency & energy consumption	62
Figure 4.6:	Magnetic domain-wall based interconnections	63
Figure 4.7:	Domain-wall motion interconnection design	64
Figure 4.8:	Composite free layer structures for “interference-free” domain-wall motion interconnects	65
Figure 5.1:	Schematics of MCSTD 2-bit full-adder.....	70
Figure 5.2:	Energy consumption components of MCSTD gates	71
Figure 5.3:	MCSTD based logic embedded DNA microarray application.....	75
Figure 5.4:	MCSTD gates as a “smart” image sensor for surveillance applications...75	
Figure 5.5:	(a) Atomic force microscope topographs of a nearly defect-free region (left) and highly defective region (right) in a 34 nm pitch nanowire crossbar [11] (b) schematic of CMOL FPGA.....	76
Figure 5.6:	MCSTD in crossbar array architecture.....	77
Figure 5.7:	MCSTD switching in crossbar circuit configuration.....	79
Figure 5.8:	Micromagnetic simulation of switching speed comparison.....	80
Figure 5.9:	Switching speed comparison for different material used for the free layer of the input device.....	81
Figure. 5.10:	Switching time improvement of MCSTD with perpendicular anisotropy.	81
Figure 5.11:	Operation of MCSTD based FPGA Look-Up Table (LUT).....	84
Figure 5.12:	MCSTD based reconfigurable routing fabric.....	87
Figure 5.13:	MCSTD intrinsic switching energy.....	89
Figure 5.14:	MCSTD based FPGA LUT read operation.....	91
Figure 5.15:	MCSTD based routing fabric reconfiguration process.....	92
Figure 5.16:	Voltage waveforms of MCSTD based Flip-flop.....	94
Figure 6.1:	Fabrication process flow of MCSTD gates.....	98
Figure 6.2:	MTJ film stack used for MTJ fabrication and (b) bottom contact structure (side and top view).....	100

Figure 6.3:	Fabrication issues of using negative e-beam resist for MCSTD fabrication (top) and Chrome hard mask layer as a solution (bottom).....	102
Figure 6.4:	E-beam lithography and Cr hard mask layer preparation for ion-milling	103
Figure 6.5:	Ion-milling process and shadow area.....	104
Figure 6.6:	MCSTD gate designs for different logic functionalities.....	105
Figure 6.7:	SEM images of electrical contacts on the output and input MTJs of MCSTD gates.....	106
Figure 7.1:	MCSTD measurement sequence and measurement parameters.....	109
Figure 7.2:	Measurement setups for MCSTD gate switching.....	110
Figure 7.3:	Mechanism of switching voltage shift in MCSTD.....	112
Figure 7.4:	Input signal dependence of switching voltage.....	114
Figure 7.5:	Input signal dependence of switching voltage.....	115
Figure 7.6:	Voltage shift of MCSTD gate versus H_{noise}	116
Figure 7.7:	Input dependent switching voltage shift vs. H_c	118
Figure 7.8:	Energy barrier height of MCSTD gates extracted from J_c vs. pulse width relations.....	119
Figure 7.9:	Input signal dependent switching speed difference.....	120
Figure 7.10:	Experimental demonstration of switching probability change versus the input device magnetization directions.....	121
Figure 7.11:	Proposed voltage operation range of MCSTD gates.....	123
Figure 7.12:	Experimentally measured time response of NAND MCSTD gate	124
Figure 7.13:	Experimentally measured time response of NOR MCSTD gate.....	125
Figure 7.14:	Total energy consumption of MCSTD gate calculated from experimental data in Fig. 7.9 and Energy consumption projection of MCSTD gate ...	127
Figure 7.15:	Device dimension of MCSTD gate at super-paramagnetic limit.....	128
Figure 8.1:	Geometrically relaxed partially unzipped Carbon nanotube (CNT).....	131
Figure 8.2:	Cross-sectional schematic of simulated GNR/CNT heterostructures to study the effect of GNR and CNTs on the tunneling FET performance	132

Figure 8.3:	Flow chart of multi-scale simulations for heterojunction Tunneling FETs	132
Figure 8.4:	Symmetric energy bands in Homojunction TFETs Tunneling happens for both $V_{gs} > 0$ and $V_{gs} < 0$	133
Figure 8.5:	Band diagram (white solid lines) and local density of states of GNR/CNT tunneling FETs	134
Figure 8.6:	Comparison of I-V curve of Carbon based tunneling FETs.....	135
Figure 8.7:	Ambipolar I-V characteristics of homojunction TFETs	135
Figure 8.8:	CNT-GNR-CNT heterojunction TFET I-V characteristics.....	136
Figure 8.9:	I-V comparison between completely unrolled GNR (“flat”) and rolled GNR.....	136
Figure 8.10:	Comparison between single (GNR-CNT-CNT) and double heterojunction (GNR-CNT-GNR) TFETs.....	137
Figure 8.11:	V_{ds} (left) and channel length (right) dependence of GNR/CNT Heterojunction TFETs.....	137
Figure 9.1:	ERSA hardware architecture.....	144
Figure 9.2:	ERSA computation model.....	146
Figure 9.3:	ERSA hardware prototype.....	150
Figure 9.4:	ERSA computation accuracy.....	151
Figure 9.5:	Output images of Bayesian Network Inference with (a) 100% (b) 90% accuracy.....	152
Figure 9.6:	ERSA execution times (a, c, e) and the corresponding zoomed-in plots (b, d, f).....	153
Figure 9.7:	ERSA L1 data cache organization (a), data cache error experiment results. (b,c)	154
Figure 10.1:	Sensory swarm. Trillions of simple devices spread in the environment adding sense to the internet.....	162

Chapter 1.

Introduction

Over the past few years, CMOS microprocessor performance improvement has slowed down. Figure 1.1 shows the performance benchmark result of all CMOS microprocessors built since 1988 [1]. Although Moore's law was about the density of transistors, it also proved to be true as in this plot that microprocessor performance doubles every 18 month or so. However, in recent years, there is a clear trend that the performance scaling of uniprocessor is saturating.

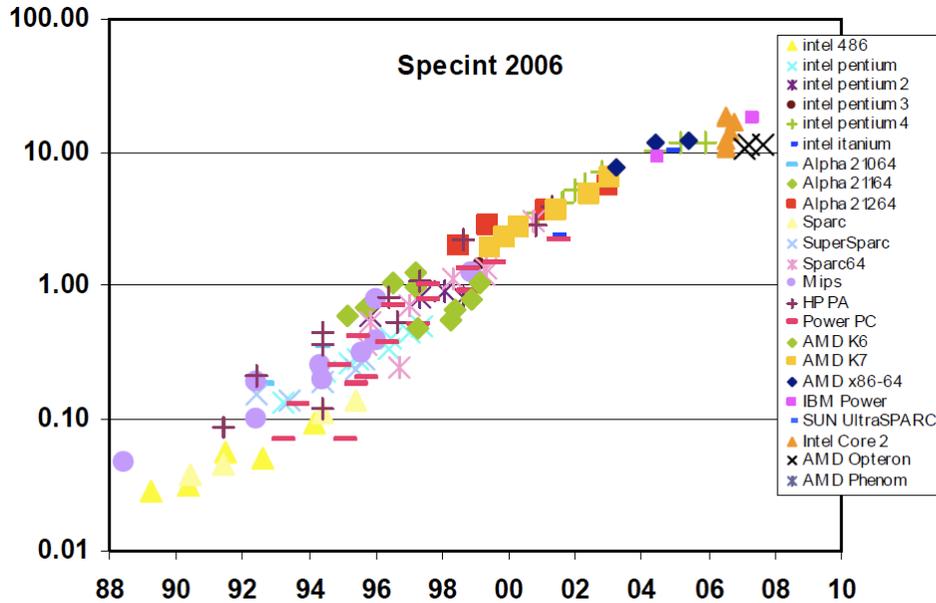


Figure 1.1. Specint2006 (www.spec.org/cpu2006/results) benchmark results of modern CMOS microprocessors [1]

The reason this occurred is because we are living in the era of power limited scaling. In the simple concept of constant electric field scaling for MOS transistors [2], all of the physical dimensions are reduced by the same amount while the body doping is increased and the applied voltage is reduced to cause the depletion regions within the devices to scale as much as the other dimensions. Unfortunately, we are now in an era where voltage is not being scaled at all for a given application [3]. Over the last ten years, it has proved to be impossible to scale V_{dd} below 1V and maintain device speed increases because of constraints on the threshold voltage in order to avoid raising standby power in the “off” transistors. This is a fundamental issue for MOSFET scaling because the thermal energy, kT/q is not the kind of thing that scales. Hence, the subthreshold slope (SS) is limited to be greater than 60mV/dec and if one tries to scale V_{th} , leakage current

increases exponentially. We have already come to the point where V_{th} optimally balances leakage and dynamic energy and V_{th} can no longer be scaled.

As a result, when we scale devices by the scaling factor k , we find circuit power **constant** with scaling, power density rising as k^2 , and the power-delay product improving only by k . In Fig. 1.2, these metrics are plotted versus gate length. The growth rate of performance (clock frequency) per power dropped from a (classical scaling) $\sim L^3$ to below $\sim L^2$ as we enter the 65-nm node ($L_g \sim 35\text{nm}$) (Fig. 1.2(a)) [3]. The improvement in performance per power density not only slowed but, actually reversed to become degraded (Fig. 1.2(b)). Power-delay product or energy per operation is reported to have actually increased after 32nm (Fig. 1.2(c)) [4]. Another reason for power degradation is the super frequency scaling: the clock frequency of microprocessors was scaling faster (by deep pipelining, etc) than what device feature size scaling allowed (Fig. 1.2(d)). Currently, device size scaling still continues but it is mainly because the cost per device continues to scale down.

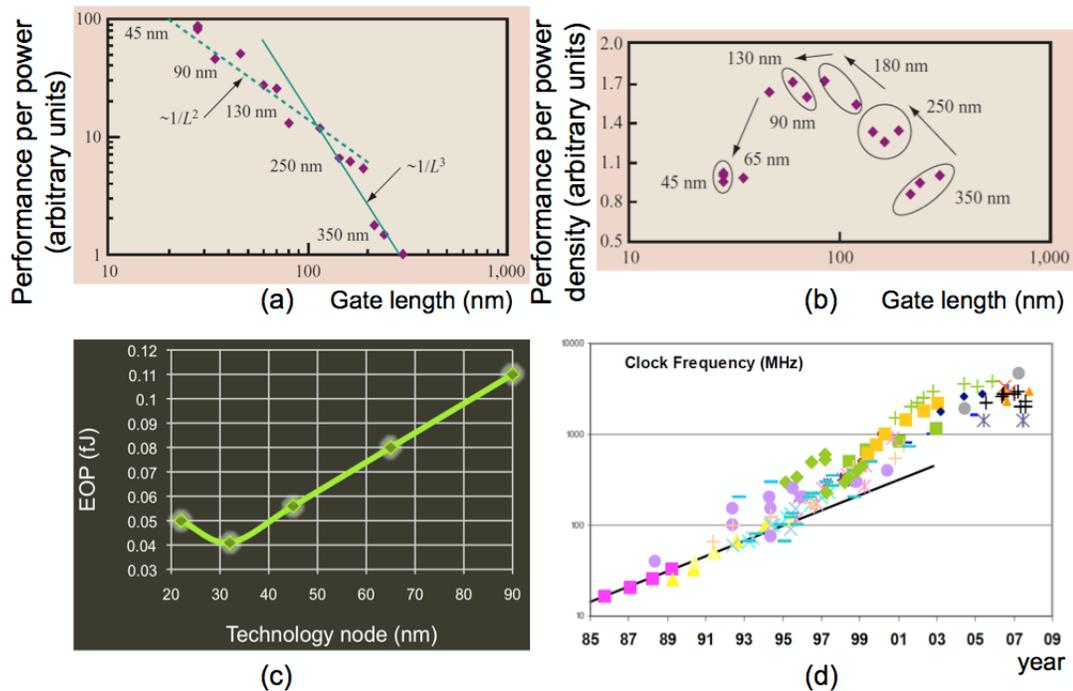


Figure 1.2 Power related metrics for MOSFET scaling. (a) Performance (frequency) per power per circuit (b) Performance per power density [3] (c) Energy per Operation (EOP) versus CMOS technology node based on actual and predictive models [4] (d) Clock frequency scaling of CMOS microprocessors [1]

In summary, CMOS uniprocessor performance is not improving because of the fixed V_{dd} and V_{th} , which cause the power density to increase despite the device scaling. At this power density, super clock frequency scaling can no longer be sustained for performance improvement. Therefore, in order to improve the CMOS performance, we have to tackle the fundamental issues of CMOS. Possible approaches are

1. **Change the fundamental design** of the MOSFET to lower the SS below 60 mV/dec and continue to scale V_{th} and V_{dd} down
2. **Scale down V_{dd}** in the current MOSFET design and deal with the consequences, e.g., timing violation errors, etc
3. **Invent a completely new device** architecture that has new state variables, such as electron spins, so that it avoids these fundamental MOSFET issues

This thesis investigates all three possible approaches to understand wide aspects of low-power research. In detail, three approaches are

1. Magnetic Coupled Spin-Torque Device (MCSTD) [7], which is a spin-based non-volatile logic whose energy barrier is manipulated by magnetic coupling between three spin-torque devices (Chap. 2~7).
2. Graphene Nano-ribbon (GNR) / Carbon Nanotube (CNT) heterojunction based Tunneling FET (TFET) [5], which achieves $SS < 60\text{mV/dec}$ via current tunneling through type-II heterojunction at the partially unzipped CNT (Chap. 8).
3. Error Resilient System Architecture (ERSA) [6] that can aggressively scale down supply voltage by masking out induced errors with cross-layer optimizations, i.e., error-resilient algorithm in probabilistic applications and asymmetric reliability multi-core hardware architecture (Chap. 9).

A wide range of “out-of-the-box” thinking is important for low-power research for the following reasons. First, conventional low power research has been mainly about how to reduce wasted power. This includes techniques such as *energy-delay trade-offs*, *multiple supply voltage levels* and *power gating*, etc [8]. However, now the wasted power

has been largely eliminated and these techniques are running out of steam. For further power reduction, fundamentally different measures have to be taken. Second, our society is changing from “technology-driven society” towards “market-driven society”, which means that important features in consumer products will be decided upon by customer’s needs rather than what is considered as the most “technologically advanced”. For example, most of recent computing devices have virtually saturated processor performance but continues to sell because other features such as memory and storage capacity, battery lifetime, user-interface, marketplace for software applications, etc. are still growing. This trend is more distinct in the domain of hand-held mobile devices. For instance, the market dominance of Apple iPhone (as of 2010) is not coming from its CPU clock frequency but from its graphical-user-interface and on-line application store.

This trend calls for low-power techniques that are tailored for “smart” devices. Future electronic systems will be dominated by the convergence between different components for better functionality and power reduction. One example is the universal memory concept that combines different memories in a chipset to reduce cost and power [9]. Research on Magnetic Coupled Spin-Torque Devices (MCSTDs) (Chap. 2~7) takes this idea one step further and provides a solution how “memory and logic” or “memory, logic and sensor” can be combined to save communication delay and power. Furthermore, new low-power techniques should benefit not only microprocessors but also a wide range of components at the system level. In this regard, the Error Resilient System Architecture work (in Chap.9) looks at voltage over-scaling and error resilience of both logic and cache memory of a microprocessor.

This thesis investigates future low-power research with the emphasis placed on the spin-based nonvolatile logic device work and its device modeling and fabrication results.

References

1. M. Horowitz, "Why Design Must Change: Rethinking Digital Design," *1st Berkeley Symposium on Energy Efficient Electronic Systems* (2009)
2. R. Dennard et al., "Design of Ion Implanted MOSFETs with Very Small Physical Dimensions," *IEEE J. Sol. St. Circ.* Vol. SC-9, pp. 256-268 (1974)
3. W. Haensch et al., "Silicon CMOS devices beyond scaling," *IBM J. RES & DEV.* Vol. 50 No. 4/5 (2006)
4. J. Rabaey, "Statistical Computing, the alternative road to load energy," *Proc. Design Automation Conference* (2009)
5. L. Leem, A. Srivastava, S. Li, G. Iannaccone, J. S. Harris, G. Fiori, "Multi-scale Simulations of Partially Unzipped CNT Hetero-junction Tunneling Field Effect Transistor," *Proc. Inter. Electron Device Meeting* 2010 (accepted)
6. L. Leem, H. Cho, J. Bau, Q. A. Jacobson, S. Mitra, "ERSA: Error Resilient System Architecture for Probabilistic Applications," *Proc. Design Automation and Test in Europe* (2010)
7. L. Leem, J.S. Harris, "Magnetic Coupled Spin-Torque Devices and Magnetic Ring Oscillator," *Proc. Inter. Electron Device Meeting* DOI: 10.1109 (2008)
8. J. Rabaey, "Low Power Design Essentials," *Springer*, DOI 10.1007 (2009)
9. S. Kang, "Embedded STT-MRAM for Mobile Applications: Enabling Advanced Chip Architectures," *Non-volatile Memories Workshop* (2010)

Chapter 2.

Introduction to Spintronics and Non-volatile Applications

2.1 Motivation

This chapter briefly reviews some of the basic principles of Spintronics. It discusses the mechanisms and the benefits of spintronic devices, which could be useful in understanding the topics of Chap. 2~4, Magnetic Coupled Spin-Torque Devices (MCSTDs), are a new logic device architecture with an alternative state variable, spin. To begin with, let's understand how computation is done in conventional CMOS circuits. The computations in CMOS integrated circuits (ICs) use the movement of charges under the influence of electric fields [1]: device operation mechanisms are based on how to block or unblock these charge movements. Consequently, device performance depends on the effectiveness of the barrier faced by the carriers and the propagation time of charges through the device. Although, device material and specific methods of controlling the barrier can vary, charge based devices will have similar basic device operation mechanisms. In other words, as long as charge movement is used to represent information, any new device idea can only make incremental improvements over CMOS.

Fundamental change in the device architecture can be realized by identifying a new computational state variable other than the electronic charge. One of the most promising candidates for this alternative state variable is the *spin* or *magnetic moment*. Spin angular momentum, or simply spin is another intrinsic property of the electron in addition to its mass and charge. There are peculiar attributes of spin that make it an attractive candidate to supplant charge-based devices. First, spin angular momentum can be transferred between particles. It allows a spin signal to be transferred from a point A to B without transporting the particle. This opens up the possibility of building a logic device without

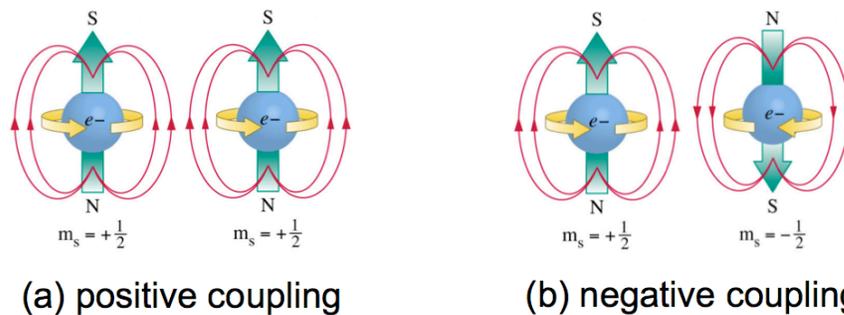


Figure 2.1 Exchange coupling of electron spins (a) positive coupling (b) negative coupling

movement of charge, which possibly can save energy, given that the energy to generate the spin signal does not offset the energy saving. Second, neighboring spins are coupled either positively or negatively. Depending on the material property, there is a tendency to align the neighboring spins in the same direction or the opposite direction. This is called *exchange coupling*, whose property extends to the macroscopic level and bi-layer ferromagnetic devices manifest two stable resistance states called magneto-resistance (MR) when two magnetic layers are positively coupled or negatively coupled (Fig. 2.1) (MR can vary as $\cos(\theta_1 - \theta_2)$, (θ_1, θ_2 : magnetization angles) .When MR shows a hysteresis, the device becomes non-volatile.

Ferromagnetic logic is not a new idea; it was commercialized over 60 years ago. A magnetic logic called Multi-Aperture Device (MAD) (Fig. 2.2) [2] was preferred for critical systems such as Canadian National Railroad Hump Yard, the New York subway and the B-50 friend or foe system [1]. One of the limitations of these magnetic logic devices is that they are not scalable. First, electron spins or magnetizations are usually less coherent than electron charges and accurately controlling spin has been a relatively more difficult task than transporting electrons for a long time. Second, controlling them with an external magnetic field inside a chip is not scalable because, a) generating magnetic fields with current in a wire takes up large space b) as the device packing density increases, it takes larger power to accurately address a bit.

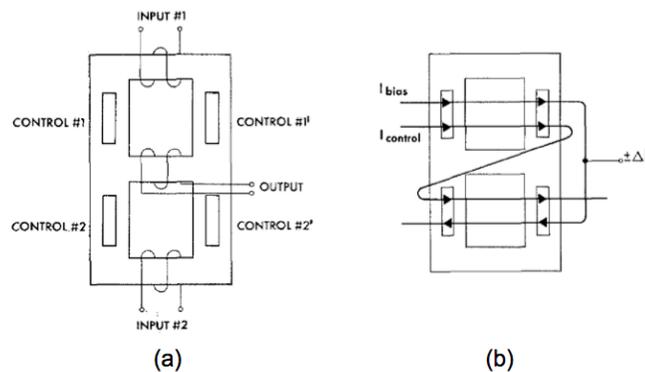
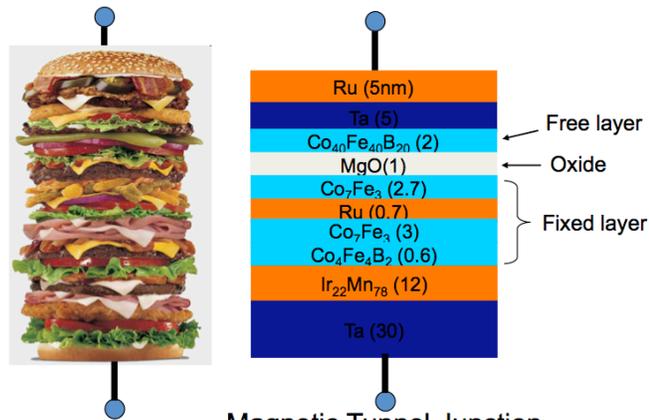


Figure 2.2 Typical Multi-Aperture Device [2]. An ac signal is applied simultaneously to both input windings. The signal is applied in such a phase that two input fluxes tend to cancel each other in the center branch. Therefore, no signal flux will be sensed in the output coil in case the bridge is balanced. If the bridge is unbalanced, some signal from one input will be transmitted to the output. Thus, the device can be used for amplitude control of a high-frequency signal

After 40~50 years, the situation has now changed for magnetic logic devices. With the advent of nanofabrication technology, it is easy to fabricate devices that can contain only a single magnetic domain, e.g., nano-pillar structure with 100~200nm dimensions. Furthermore, magnetic devices can be electrically controlled with a phenomenon called *spin-torque transfer* [3, 4], which will be discussed in detail in later sections. As a result, scalable nanomagnetic logic devices including those termed “Spintronic” devices are finally achievable by exploiting the extensive CMOS fabrication processes that are available today. For example, Fig. 2.3 shows one of the most successful spintronic devices called a Magnetic Tunnel Junction (MTJ). It is a two terminal device with a sandwich structure of two ferromagnetic metals like Co, Fe, Ni that are separated by a tunneling oxide in the middle. It has many layers, perhaps as many as the burger on the left. This device is interesting because,

1. It has hysteresis between two resistance states
2. It is a magnetic device, whose magnetization can be controlled electrically

Next two sections explain the origins of these properties.



Magnetic Tunnel Junction
(figure inspired by a slide from Intermag '09)

Figure 2.3 Magnetic Tunnel Junction. Two-terminal spintronics device widely adopted for MRAM and hard disk drive readhead applications. It consists of two ferromagnetic layers (free, fixed layers), tunneling oxide layer and many others.

2.2 Spin-orbit interaction

How can we make use of electron spins as an electrical signal? One possible way is to use spin-dependent characteristics, such as, *spin-dependent scattering*. One of the underlying physical mechanisms behind spin-dependent scattering is called “spin-orbit interaction or coupling” [5]. Spin-orbit interaction describes the effects of an electron’s orbital motion on the orientation of its spin. From the electron rest frame, it “sees” a positive charge in motion just as the sun appears to be in motion about a casual observer on the earth. Hence, the electron is situated near the center of a current loop, which generates a magnetic field that causes a preferred direction of orientation for the spin magnetic moment of the electron. An induced magnetic field thus acts differently on spin-up and spin-down electrons.

When there is an electric potential induced by impurity charges, spin-up and spin-down electrons scatter in different directions (Fig. 2.4). Therefore, the length for which the spin of an electron is conserved is finite. This length is called the *spin-flip mean free path* and typically takes values in the range 100 nm ~ 10 μm. Due to scattering of electrons, the length an electron travels with a fixed spin direction is much shorter than the spin-flip mean free path. This length is called *spin-diffusion length*, λ_{spin} . To find the spin-polarized current in a non-magnetic metal, it is necessary that the system length, L , be much shorter than λ_{spin} [6].

Spin up/down electrons can be separated by spin dependent scattering. This

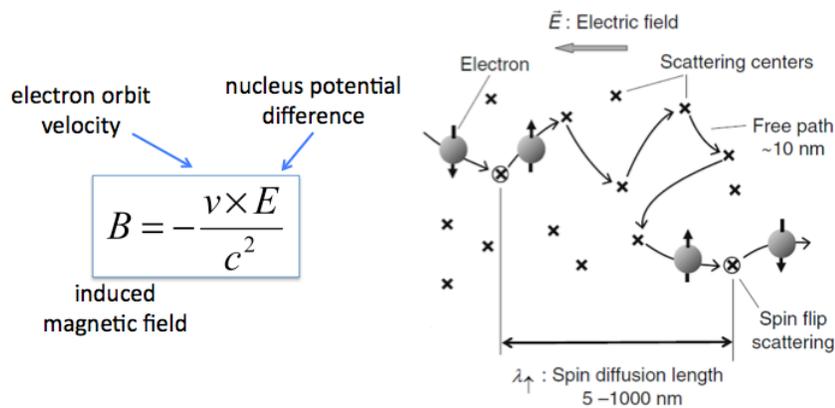


Figure 2.4 Spin-Orbit Interaction induced magnetic fields. During scattering process, the spin of electrons can get flipped or electrons can get scattered to different directions depending on the spin

phenomenon is called the *Spin-Hall effect* [7,8,9] which has some resemblance to the original Hall effect. In the classical Hall effect experiment for charge current, the orthogonal electric and magnetic fields induce Lorentz forces on positive and negative charges, which get accumulated on opposite surfaces of the sample. Similar phenomenon happens with spin up and spin down electrons. Here, the magnetic field is not necessary because effective magnetic fields are induced from the spin-orbit interaction. On the contrary, if a strong enough magnetic field is applied in a direction perpendicular to the orientation of the spins at the surfaces, the Spin Hall Effect will disappear because of the *spin precession* around the direction of the magnetic field [8,9]. Due to the spin-orbit interaction, which leads to the coupling of spin and charge currents, an electrical current induces a transverse spin current (a flow of spins) and vice versa [7,8,9] (Fig. 2.5). One can intuitively understand this effect by using the analogy between an electron and a spinning tennis ball, which deviates from its straight path in air in a direction depending on the sense of rotation (the Magnus effect) [10]. However, no Hall voltage can be measured between these opposite spin electrons. This is because the electrons are all negative charges, no potential can be measured between the two points where spin polarized electrons are accumulated. Thus, the Spin-Hall effect is quantified indirectly by measuring the Kerr rotation in reflected polarized light. Spin-Hall effect can be used to generate “dissipation-less” spin current. There is no net current flow because two spin types of electrons flow in the opposite direction. With the Spin-Hall effect, spin angular momentum can be transferred and current induced magnetization reversal is possible.

In summary, any disorder at interfaces or random distribution of electrode and oxide

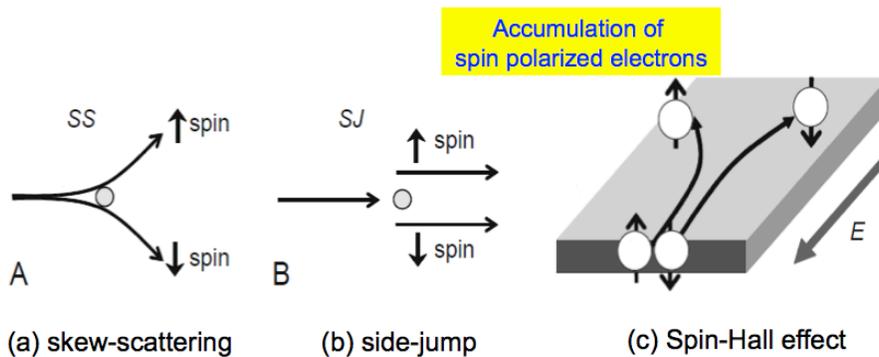


Figure 2.5 Spin-dependent scattering and Spin-Hall effect. At impurity induced electric potential, spin-up/down electrons scatter in the opposite directions (a, b). Spin-Hall effect generates “dissipation-less” spin current in normal direction to charge currents (c)

atoms at the interface work as a disturbance to potential. Spin dependent scattering is the outcome of spin orbit interaction at these non-uniform interfaces.

2.3 Spin dependent electronic structures

Another reason for spin-dependent scattering is that the electronic band structure can be different for up-spin and down-spin electrons. Figure 2.6 shows density-of-state (DOS) versus energy plots. Parabolic plots are the DOS of 4s orbital and half ellipsoids are those for 3d orbitals. In most materials up-spin and down-spin electron DOS are symmetric: there are always equal number of up-spin and down-spin electrons. But the situation is different in magnetic materials. Density of states (DOS) is asymmetric for magnetic materials: electron bands of up-spin electrons are at lower energy compared to the down-spin electrons. This means there are many more up-spin electrons than the down-spin electrons. Therefore we have net magnetization. [Note: '+' means majority spin not the 'spin up' and '-' is for minority spin electrons. Please don't confuse majority spin with spin-up electrons. Majority spin can be either spin up and down electrons and it changes whenever the magnetization of a material changes.]

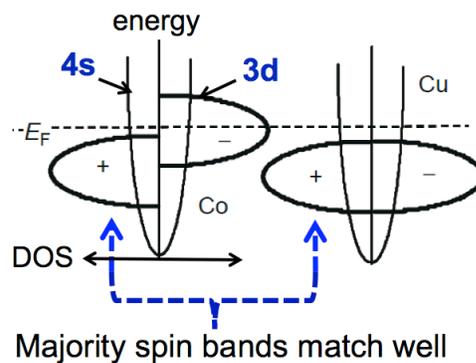


Figure 2.6 Density of States (DOSs) comparison of magnetic (Co) and non-magnetic (Cu) materials. For magnetic material, 3d orbital DOS is asymmetric. When Co and Cu make interface, majority spin bands match well but, minority spin electrons don't, which results in spin-dependent scattering [6]

Consider, for example, Co and Cu layers interfaced with each other or Cu layers having Co impurities, the energy bands of one type of spin electrons will match but the other type will have a mismatch in the electronic states. In this case, the majority spin electrons have matching energy bands but the minority spins don't. As a result, the

impedance is low for majority spins and high for minority spins, which leads to spin-dependent scattering.

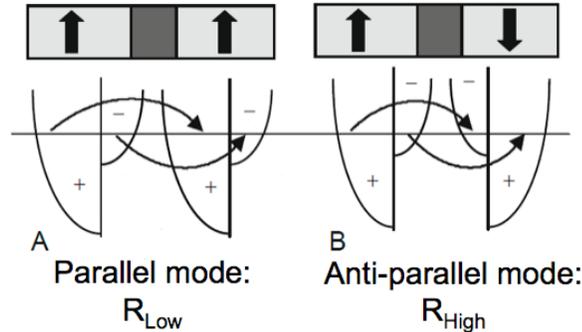


Figure 2.7 Density of States (DOSs) comparison of magnetic tunnel junctions. In parallel mode, one spin-type electron is majority in the both electrodes, which leads to low overall resistance. In anti-parallel mode, both spin-type electrons are majority in one side and minority in the other, resulting in high total resistance [6]

When one has two ferromagnetic layers separated by an oxide layer, the spin-dependent scattering will depend on the magnetization alignments of the two ferromagnetic electrodes (Fig. 2.7). Tunneling current is proportional to the product of the density of states in the two electrodes. If two ferromagnetic layers are in the same spin orientations (*parallel mode*), the majority type of spin in one electrode will be a majority in the other making the product of the density of states for this type of spin electrons large – overall resistance is low. In contrast, if two ferromagnets are in the opposite spin orientations (*anti-parallel mode*), both type of spins will be majority in one electrode but minority in the other. The product of the DOS for spin up and spin down electrons are equally small – overall resistance is high. This change in resistance due to the difference in magnetization alignment is called “magneto-resistance”. Magneto-resistance (MR) is an essential element in almost all spin devices: spin sensors, memories, even logic devices depend on the MR for readout. In this discussion, two spin-type electrons are treated independently. For this treatment to be valid, the basic assumption is that the spin flip mean free path \gg electron travel distance.

At high temperature, both spin energy bands will be filled. This results in an equal number of spin-up and down electrons and so, the material will no longer be spin-

polarized or magnetized. This temperature is called the *Curie* temperature (T_C) for ferromagnetic materials and the *Neel* temperature (T_N) for anti-ferromagnetic materials.

2.4 Magneto-resistance (MR)

Ferromagnetic structures of spin-torque devices that can switch between negative/positively-coupled configurations have been pursued for larger MR signal and faster switching speed. In the early days, one method used by researchers was the interlayer (usually anti-ferromagnetic) coupling (Fig. 2.8). Two layers are negatively aligned when no magnetic field is applied. When a moderate external field is applied, two layers tend to align in one direction. But, interlayer coupling makes the free layer hard to switch and the switching speed was slow.

Another method was to use the difference in the magnetic coercivity of the material. In other words, magnetically soft and hard materials are used for the free and fixed layers. At low or medium field, the soft magnetic material will switch first, which results in a negatively coupled configuration. At high fields, the hard magnetic layer switches and the two layers become aligned. In order to avoid the interlayer coupling, a non-magnetic layer is inserted between the ferromagnetic layers. This structure is called a “spin-valve”. With the introduction of spin-valves, the observed magnetoresistance increased significantly from 1% to 20% at room temperature. This unusually large increment was thus given the name “giant magneto-resistance (GMR)”. Spin-valves are one of the most successful spintronic devices and are used for GMR sensor. Comparing the scaling of microprocessors and hard-disk drives (HDDs), CPU performance has saturated over recent years while HDD density has not (Fig. 2.9) [12]. The enabling technology behind the successful scaling of HDDs is the spin-valve. It is used in HDD readheads, the magnetic sensor that senses stray magnetic signals from bits in HDDs. The minimum bit pitch size significantly decreased by shifting HDD technology from coil-based readheads to spin-valve readheads, which can be a few hundred nm wide [12].

Magnetic Tunnel Junctions (MTJs) make even more attractive sensors than spin-valves for high-density readheads because of the larger signal from a greater tunnel-magneto-resistance (TMR). State-of-the-art TMR ratio is over 600% while GMR is still less than 100% [12]. In addition, MTJs are better suited for Current-Perpendicular-to-

Plane (CPP) configuration when spin-valves are usually used in Current-In-Plane (CIP) (Fig. 2.10). The CPP head has a narrower shield spacing, which is a significant advantage over CIP devices. On the other hand, the large MTJ sensor resistance is accompanied by large noise, including shot noise, which is specific to tunneling sensors. In addition, the high RC time constant of a high-impedance TMR sensor results in an unacceptable decrease in attainable sensor bandwidth. Therefore, the signal-to-noise ratio (SNR) for high-density MTJ head sensors is sufficiently high for head applications only when the resistance-area (RA) product is sufficiently small ($<10 \Omega\text{-}\mu\text{m}^2$). Thus, only low-resistance MTJs are suitable for recording heads.

TMR is larger than GMR because of the additional spin-filtering effects

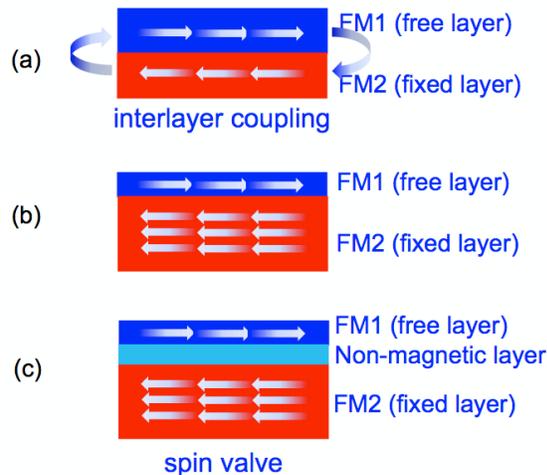


Figure 2.8 Evolution of magnetoresistive film structures. (a) interlayer coupling based (b) coercivity difference based (c) spin-valve

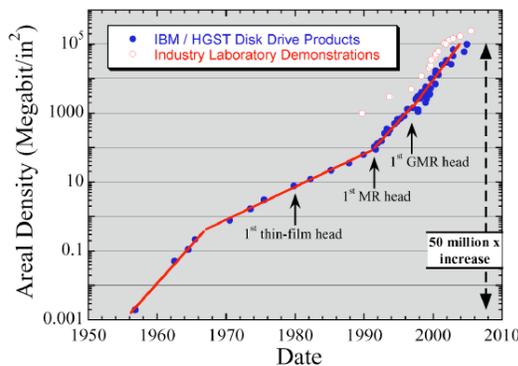


Figure 2.9 Areal-density growth curve of Hard Disk Drive (HDD) recording products [12]

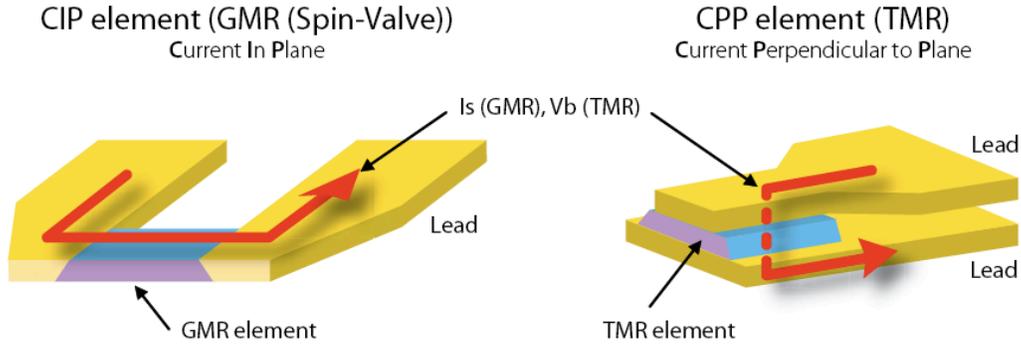


Figure 2.10 Current-in-plane (CIP) vs. Current-perpendicular-to-plane (CPP) readhead architecture [13]

sorinating from tunneling. First, the electronic states of ferromagnets depend on spin, the tunneling probability becomes spin dependent due to this momentum filter effect, and the decay rate may be written as

$$\kappa_s = \sqrt{2m(U - E_s(k_z))} / \hbar$$

The spin dependent κ_s can produce high TMR ratios.

Second, if we assume the component of the wave vector k_{\parallel} parallel to layer planes is conserved in tunneling (a.k.a *specular tunneling*), only states on the Fermi surface with the same k_{\parallel} in the free and fixed layer may contribute to tunneling. In anti-parallel (AP) alignment, the state with the same k_{\parallel} may be included in the one spin state, but not in the other spin state. In this case the tunnel conductance becomes small and TMR becomes large. Third, the difference between the symmetries of the wave functions of tunneling electrons in the free and fixed layers exerts a strong influence on TMR. The wave function of each state in metals has a specific symmetry and electrons with a certain symmetry are not able to transfer into a state with different symmetry. Therefore, electrons on the free layer can tunnel through the barrier into the fixed layer only when the states specified by k_{\parallel} in the free and fixed layers have the same symmetry. For AP alignment, the states on the Fermi surfaces of the free and fixed layers can have different symmetries, e.g., the symmetry in wave function exists in the majority spin band but does not appear in the minority spin band, which leads to low tunneling conductance [6]. Lastly, TMR is also related to the strength of the local chemical bonds formed at the interface with the tunnel barrier. For Co-Pt ferromagnetic alloy, Co is magnetic but Pt is

not. Since Pt forms a much weaker bond with oxygen than does Co, we conclude that the tunneling current from Co-Pt alloys with alumina tunnel barriers will be dominated by tunneling from the Co. On the other hand, Co-V alloy with AlN as a tunnel barrier, V (vanadium) forms a much stronger bond with O than does Co, so the tunneling current from V would increase rapidly, which reduces the spin polarization of net tunnel current [14].

To date, the MgO tunneling barrier has been experimentally reported to yield the highest TMR ratio compared to any other materials. Compared with other oxide barriers (e.g., Al₂O₃) the characteristics of MgO can lead us to insights on what are the specific requirements to obtain high TMR. These peculiar characteristics are:

1. Ionic bond, e.g., MgO has ionic bonds while Al₂O₃ has covalent bonds
2. Fewer impurities at oxide interfaces. Any impurities work as spin-flipping centers
3. Crystalline structure (MgO) is preferred over amorphous (Al₂O₃)

Together with the electronic structure of the ferromagnetic electrode that determines the magnitude of the spin dependent tunneling, the spin polarization of the tunneling current can be strongly modified by tunneling matrix elements that themselves depend on chemical bonding at the ferromagnet/tunnel barrier interface and also on the symmetry of the electrode's conduction band wave functions. Weakly magnetic metals can give rise to highly spin-polarized currents by suitable wave function or chemical bond engineering. Similarly, strongly magnetic metals may be tuned to result in only very weakly spin polarized tunnel currents [14].

2.5 Spin-torque transfer effect and magnetic reversal

As mentioned earlier in Sec. 2.1, the spins of neighboring electrons are “*exchange*” coupled : they have a tendency to align the neighboring spins in the same direction or the opposite direction. Now let's extend our discussion to the case when spin-polarized electrons are injected into non-polarized electrons. As the electrons from two groups encounter, non-polarized electrons increasingly become spin-polarized into the same spin directions as the injected electrons due to exchange coupling. This phenomenon can be

macroscopically understood that there is a “*transfer*” of spin angular momentum and this time change in the spin is called *spin-torque*. Spin-torque transfer effect has a huge theoretical and practical significance, because it opens up a way to electrically control magnetic devices.

Spin-valves and MTJs are the most successful spin-torque devices. They are all in nano-pillar shape. Nano-pillars (cross-section of $100 \times 100 \text{nm}^2$) are suitable structure to observe spin-torque effect because,

- i) it has a single magnetic domain, which is less random and has controllable magnetization
- ii) it increases current density.
- iii) spin-torque transfer phenomenon will be more effective at smaller dimension

Magnetic devices with an elliptical shape have an associated energy called *magnetic shape anisotropy*. Magnetizations in all the sub-domains tend to align to the longer axis. This can be explained by considering surface magnetic charges – imaginary magnetic charges always positioned at the surface or boundaries of the device. If the magnetization is aligned along the short axis, (i.e., the hard-axis) 1) the distance between the magnetization and the surface charge is short and 2) surface area also larger – more “magnetic” charges. These result in high magnetic potential energy. If the magnetization is aligned along the longer axis, (i.e., the easy axis), distance to the surface charges is large and the magnetic energy becomes smaller and magnetizations aligned along the easy-axis becomes more stable than in the hard-axis case. This magnetic shape anisotropy generates an energy barrier that must be overcome for magnetic reversal (Fig. 2.11). Materials such as CoFeB - the commonly used free and fixed ferromagnetic layer

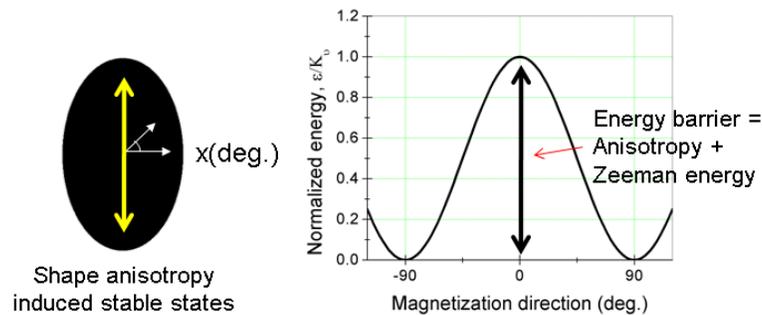


Figure 2.11 **Energy barrier of nano-pillar shaped spin-torque device.** The elliptical shape induces two stable states at $+90^\circ$ and -90° magnetizations

material for spin-torque devices, only has magnetic shape anisotropy. But other materials have an additional *magnetic crystalline anisotropy*; in which the magnetization favors a particular crystal orientation. In this case, magnetic crystalline anisotropy energy adds to the energy barrier height. In addition, external magnetic fields alter the energy barrier height and the potential corresponding to these external fields is added to it.

For magnetic field driven switching, energy required for magnetic reversal equals the energy barrier height (plus/minus thermal fluctuation). For spin-torque transfer based switching, the efficiency of spin-torque transfer has to be included. The magnetic reversal process can be considered this way. Electron spins of the magnetic layer to be switched are loosely coupled. Each of them is coupled to neighboring spins with exchange coupling but thermal fluctuations affect them individually. In addition, spin-torque transfer or spin angular momentum transfer from incoming spin polarized current occurs to individual electrons. This makes the magnetic reversal process rather incoherent and random switching. As an analogy, magnetic reversal with spin-torque transfer is like herding a group of flies that are moving in a group but flying on their own. Finally, switching is the interaction between applied energy versus restoration force. Restoration or damping force tends to bring the magnetization back to the initial state. This is why sufficient energy in terms of spin angular momentum has to be provided before switching happens. When a voltage pulse is used to switch spin-torque devices, driving voltage and pulse width are in an inverse relationship. When the pulse width is increased, the switching voltage is lowered. But this is not proportional to the increased time. For example, when pulse width is doubled, the switching voltage does not become half. This is because magnetic reversal is a complex process where 1) spin angular momentum is transferred in a random process and 2) it is an interaction between the restoration force and magnetic switching process.

For simulation of nanomagnets and spin-torque transfer, the Landau-Lifshitz-Gilbert (LLG) equation is primarily used. To explain each term, first, let's consider what happens to magnetizations when a magnetic field and spin-torque are applied. Magnetization starts to precess in a direction that is normal to both magnetization and magnetic fields.

$$\frac{dM}{dt} = -M \times B$$

This rotation doesn't go forever because there is a damping. Landau and Lifshitz's idea was that damping works to align M to B so, they added another term normal to $M \times B$.

$$\frac{dM}{dt} = -M \times B - M \times (M \times B)$$

Gilbert's idea was that damping is a friction, which should be proportional to dM/dt . It turns out dM/dt equals to $M \times B$ (as shown above) and the two approaches become identical. Spin-torque can work against/for damping and keep the magnetization rotating, or it can make the magnetization switch to the opposite direction. Pictorial descriptions of each term are illustrated in Fig. 2.12. The final form of the LLG equation is shown in the equation below [15].

$$\frac{dm}{dt} = -|\gamma|m \times H + \alpha \left(m \times \frac{dm}{dt} \right) + |\gamma|\beta\varepsilon (m \times m_p \times m)$$

m = reduced magnetization, M/M_s

γ = Gilbert gyromagnetic ratio

$$\beta = \left| \frac{\lambda}{\mu_0 e} \right| \frac{J}{tM_s}$$

α = damping constant

t = thickness of the free layer

$$\varepsilon = \frac{P\Lambda^2}{(\Lambda^2 + 1) + (\Lambda^2 - 1)(m - m_p)}$$

P = polarization

If we include thermal fluctuations in the system [16],

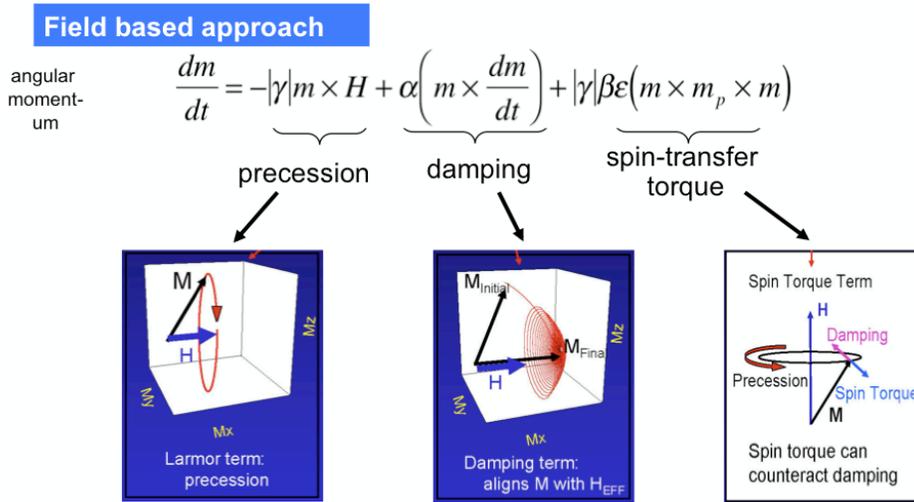
$$\frac{dm}{dt} = -|\gamma|m \times (H + H_{Thermal}(t)) + \alpha M \times (M \times (H + H_{Thermal}(t))) + |\gamma|\beta\varepsilon (m \times (m_p \times m))$$

where, thermal fluctuation is modeled as Wiener process (e.g., Brownian motion).

$$H_{Thermal}^{(i)}(t)dt = v \cdot dW^{(i)}, \quad v = \sqrt{\frac{2\alpha k_B T}{\mu_o M_S^2 V}}$$

$W^{(i)}$: Wiener process, α : damping constant

From the LLG equation, it is possible to categorize the operation modes of spin-torque



- Solve with Time evolvers : track LLG equation with time
- e.g. Euler, Runge-Kutta methods



Energy based approach

Magnetic energy

$$E = \int_V \left[\underbrace{A(\nabla m)^2}_{\text{exchange}} + \underbrace{\epsilon_K}_{\text{anisotropy}} - \underbrace{\mu_o M_S (H_{ext} \cdot m)}_{\text{Zeeman (external field)}} - \underbrace{\frac{1}{2} \mu_o M_S (H_d \cdot m)}_{\text{demagnetization}} \right] d^3r$$

- Solve with Energy minimization technique : locate local energy minima.
- e.g. conjugate gradient method

Figure 2.12. Micromagnetics modeling of magnetic reversal. Modeling can be done either by (a) Field based approach (solve Landau-Lifshitz-Gilbert equation with spin-torque term) or (b) Energy based approach (minimize total magnetic energy)

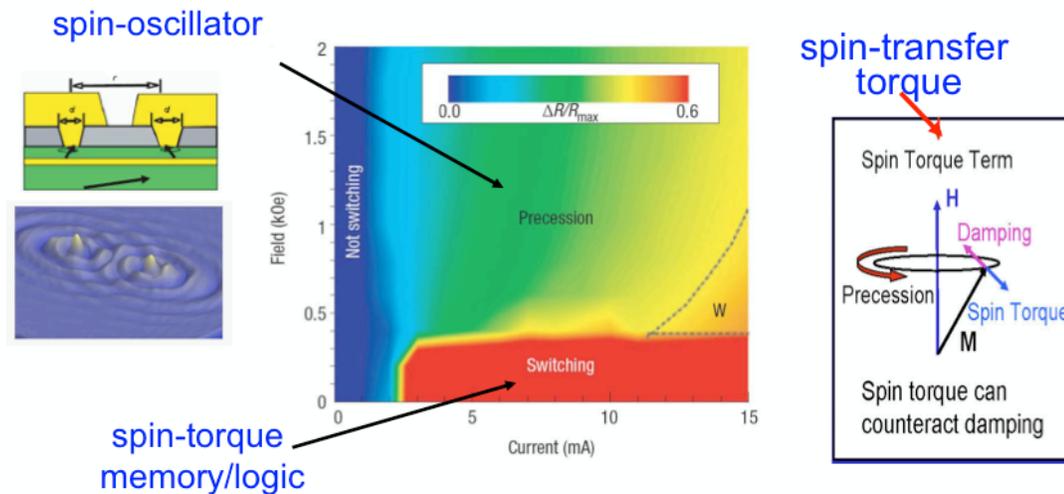


Figure 2.13. Operation modes of spin-torque devices : (a) Switching (spin torque > damping) (b) Precession (spin torque ~ damping) (c) Not switching (spin torque < damping)

devices based on the magnetic field (or damping induced from it) and spin-torque strength. If we assume the magnetization is already aligned to the magnetic field, Fig. 2.13 shows what happens to the spin-torque device when damping and spin torque compete. Magnetic precession is an interaction between magnetic damping and spin torque. As the current increases, the spin torque effect gets stronger. As the magnetic field goes up, the damping will be dominant. In the “precession” region, the two forces balance each other, which makes spin torque devices continuously oscillate. In the “switching” region, spin torque wins and switching happens, which enables a spin torque device to function as a logic or memory device.

The next example (Fig. 2.14) illustrates on how spin-torque transfer can switch a magnetization in a MTJ. The free layer magnetization switching direction depends on the direction of the current sent through this device. Let’s see how this happens. If electrons flow from the fixed to the free layer (a), the electrons are polarized in the fixed layer, and when these spin-polarized electrons enter the free layer, they exert a spin torque on the electrons already existing in the free layer. This will cause a parallel alignment of the layers. If the current is reversed (b), electrons that are initially polarized by the free layer get partially reflected by the fixed layer and return to the free layer. Since they are reflected electrons, their spin-polarizations are opposite to the transmitted electrons in (a)

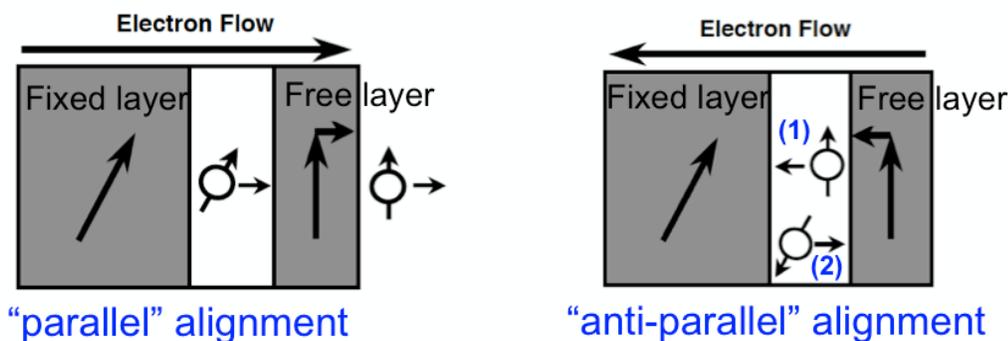


Figure 2.14. **Spin-torque transfer switching of magnetic tunnel junction** (a) electrons flowing from the fixed layer switches the free layer into parallel alignment (left). (b) electrons flowing from the free layer switches the free layer into anti-parallel alignment (right).

and the spin-torque they exert is also opposite. Now, anti-parallel alignment of the layer moments is favored by the spin-torque. As a result the magnetization and output resistance of magnetic spin-torque devices can be controlled electrically via this spin-torque transfer effect.

The spin-torque transfer effect is a relatively weaker mechanism than magnetic fields. It requires a small Resistance-Area (RA) product and low coercivity material for successful switching. Tunneling oxide needs to be thin to increase tunneling (reduce

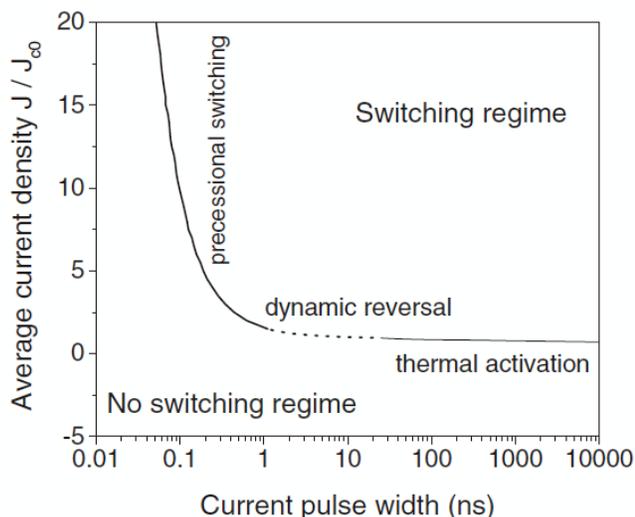


Figure 2.15. **Spin current driven magnetization switching phase diagram**. The three switching modes are thermal activation (solid line), dynamic reversal (dotted line) and precessional switching (thick solid line). The parameters are taken as $\alpha = 0.02, H_K = 500Oe, 4\pi M_S = 18kOe$ [19]

resistance) but, thin oxides also reduce the TMR ratio as well. A low switching current density is required to scale down the current driver transistors in a 1 MTJ – 1 transistor cell structure. In order to reduce the current density, a perpendicular magnetic anisotropy material has to be incorporated [17, 18].

Switching modes of spin-torque transfer can be divided into three domains: precessional, dynamic, thermal activation switching as shown in Fig. 2.15. Notice that the required current density exponentially increases as the switching time approaches sub 100ps region.

2.6 Electric field based switching

Recently research efforts are being made to develop ways to control magnetization with electric field. The motivation behind this investigation is that it could prove to be a zero or ultra low current solution compared to spin-torque transfer techniques. Multiferroic materials, such as BiFeO_3 are both ferroelectric and ferromagnetic, which

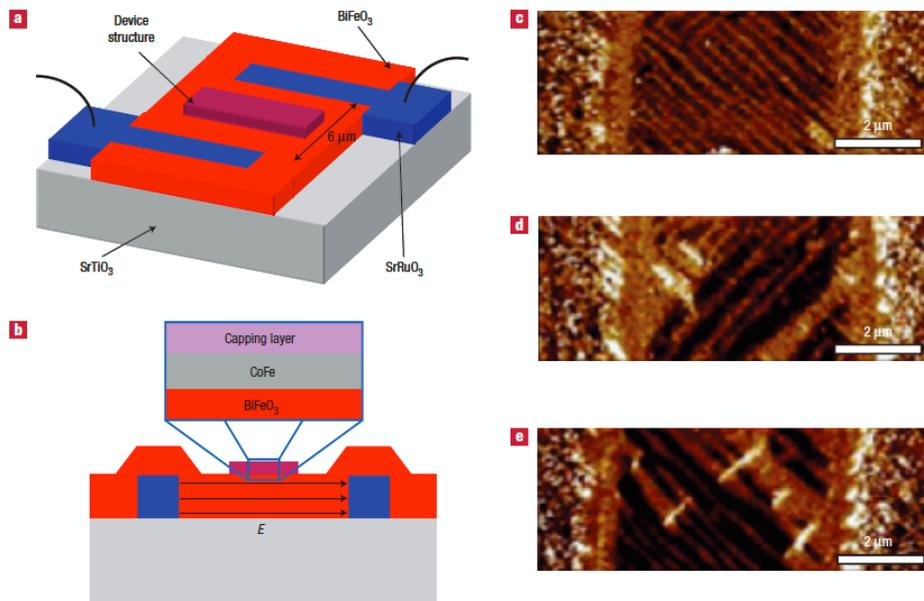


Figure 2.16. **Electric field based switching device structure.** Three-dimensional (a) and cross-sectional (b) schematic diagrams of the coplanar epitaxial electrode device showing the structure that will enable controlled ferroelectric switching and electrical control of local ferromagnetism in the CoFe features. In-plane PFM images showing the ferroelectric domain structure for a device in the as-grown state (c), after the first electrical switch (d) and after the second electrical switch (e) [20]

means that they change polarization (ferro-electric) or magnetization (ferro-magnetic) with both electric field and magnetic field. By coupling multiferroic materials with ferromagnetic ones, it is possible to make ferromagnetic materials that react to electric fields. In detail, it is possible to use two types of electromagnetic coupling phenomenon that are manifested in heterostructures consisting of a ferromagnet in intimate contact with the multiferroic material, such as BiFeO₃ (Fig. 2.16). The first is an internal, magnetoelectric coupling between antiferromagnetism and ferroelectricity in the BiFeO₃ film that leads to electric-field control of the antiferromagnetic order. The second is based on exchange interactions at the interface between a ferromagnet (Co_{0.9}Fe_{0.1}) and the anti-ferromagnet.

Experimental results reveal the possibility to locally control ferromagnetism with an electric field via [21],

1. injection or depletion of carriers that can mediate ferromagnetic alignment of the magnetic moments in oxides. Applying electric field enhances (reduces) the conductivity by injection (depletion) of carriers therefore changes the ferromagnetic alignment of magnetic materials (e.g. La_{0.9}Ba_{0.1}SRMnO₃ on Nb-doped SrTiO₃).
2. mechanical strain caused by electric field in magnetostrictive or piezomagnetic materials. Strain can be transferred to magnetic components and changes the magnetization (e.g. Ferromagnetic spinel formed by epitaxial perovskite ferroelectric matrix on bottom electrode).
3. coupling electric field induced polarization with magnetization orientation (e.g. magnetic anisotropy) through structural distortion (e.g. "spin-spiral" multiferroics, TbMnO₃).
4. inducing magnetization at the interface between a ferromagnetic metal and a nonpolar dielectric. The accumulation of spin-polarized electron adjacent to the interface with dielectric leads to net change in magnetization (e.g. SrTiO₃/LaAlO₃/La_{0.6}Sr_{0.4}MnO₃).

2.7 Spintronic logic device

As mentioned in the introduction, Spintronics originated from a study to build a logic device with an alternative state variable, called *spin*. To get a sense of how Spintronics logic devices work, let's first understand how a Spin-FET works. There exist a number of variations including the original Datta-Das spin-optical modulator [22] but, here I will introduce the one by Hall and Flatte [23].

As shown in Fig. 2.17, the source and drain are ferromagnetic materials (FM). If you send a current through a FM material, there will be spin-dependent scattering; electrons that have the same spin as the material are easy to transmit but the other spin type electrons can't be transmitted. For example, if the source region is magnetized in the spin-up direction then, it passes spin up electrons only. On the other hand, the drain region is magnetized in the opposite direction and blocks the spin up electrons. Now spin-up electrons are trapped in the channel region and no current flows and the transistor is turned OFF. When V_{gs} is applied, some of the spin-up electrons get scattered into spin-down electrons, which can pass through the drain barriers. Now the transistor is turned ON and current flows. This device can be low power because there is no energy barrier that you need to raise electrically and no subthreshold slope (SS) issues, either. High mobility current is expected in 2 DEG channel due to separated doping region and channel region.

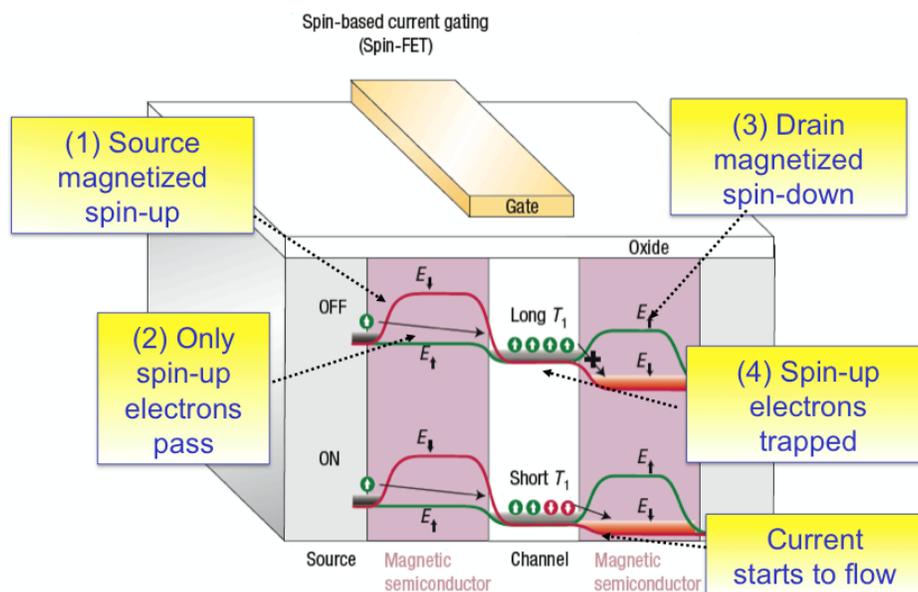


Figure 2.17. Spin-based current gating of Spin-FET [23]

Experimentally, no working spintronics transistor of this type has been fabricated over the period of twenty years of research since the time this idea was initially conceived. This is because,

1. the spin-filtering effect at source and drain is not perfect
2. spin-injection into semiconductor is still inefficient due to impedance mismatch between metal and semiconductor
3. spin information lost in highly doped channel
4. spin signals are hard to detect in the drain region. There can be high noise level because electrons that were inside the drain can also flip
5. still requires the movement electron carriers

Until these challenges are solved, there seems to be significant uncertainty as to when a working spintronics transistor might be realized. Although, continuing research on spinFET is definitely encouraged, it is also very important to conceptualize alternative solutions, such as, a new device architecture that is not based on a MOSFET type logic device architecture. One reason behind such an approach would be the fact that we still don't know what is the optimal device architecture for spintronics. In the following chapters, an alternative spintronics logic device that is free from the problems of the spinFET and can be built with today's technology, will be presented. Fabrication and measured results from a working prototype device will follow.

References

1. Broad Agency Announcement Non-Volatile Logic (NV Logic) Microsystems Technology Office, DARPA-BAA-10-42, (2010)
2. A.A. Jaecklin, "A Multiaperture Magnetic Control Device," *IEEE Transactions on Magnetics*, vol. MAG-2, no. 3 (1966)
3. J.C. Slonczewski, "Current-driven excitation of magnetic multilayers," *Journal of Magnetism and magnetic Materials*, vol. 159, issue 1-2, pp. L1-L7 (1996)
4. D.C. Ralph, M.D. Stiles, "Spin Transfer Torques," *Journal of Magnetism and magnetic Materials*, vol. 320, issue 7, pp. 1190-1216 (2008)
5. R. O'Handley, "Modern Magnetic Materials – Principles and Applications," Wiley Inter-Science (2000)
6. J. Inoue et al., "Nanomagnetism and Spintronics," DOI: 10.1016/B978-0-444-53114-8.00002-9, Elsevier B.V. (2009)
7. M. I. Dyakonov, V. I. Perel, "Possibility of orientating electron spins with current," *Sov. Phys. JETP Lett.* 13: 467, (1971)
8. M.I. Dyakonov and V.I. Perel, "Current-induced spin orientation of electrons in semiconductors," *Phys. Lett. A* 35: 459. doi:10.1016/0375-9601(71)90196-4, (1971)
9. M.I. Dyakonov, "Magnetoresistance due to edge spin accumulation" (abstract page). *Phys. Rev. Lett.* 99 (12): 126601. doi:10.1103/PhysRevLett.99.126601. PMID 17930533, (2007)
10. http://en.wikipedia.org/wiki/Spin_Hall_effect
11. N.F. Mott, Proc. R. Soc. Lond. Ser. A **153**, 699; **156**, 368 (1936)
12. J.R. Childress, R.E. Fontana, Jr. "Magnetic recording read head sensor technology," *Comptes Rendus Physique* vol. 6, issue 9, pp. 997-1012 (2005)
13. Fujitsu white paper, "CPP Read-Head Technology Enables Smaller Form Factor Storage," http://www.fujitsu.com/downloads/COMP/fcpa/hdd/ccp-based-storage_wp.pdf
14. S.S.P. Parkin, "Spin-Polarized Current in Spin Valves and Magnetic Tunnel Junctions," *MRS Bulletin*, vol. 31 (2006)

15. M.J. Donahue, D.G. Porter, "OOMMF User's Guide, Version 1.0," *Interagency Report NISTIR 6376*, National Institute of Standards and Technology (1999)
16. D.V. Berkov, "Magnetization Dynamics Including Thermal Fluctuations: Basic Phenomenology, Fast Remagnetization Processes and Transitions Over High-energy Barriers," *Handbook of Magnetism and Advanced Magnetic Materials*, vol.2: Micromagnetism, John Wiley & Sons, Ltd. ISBN: 978-0-470-02217-7 (2007)
17. S. Mangin, Y et al., "Reducing the critical current for spin-transfer switching of perpendicularly magnetized nanomagnets," *Appl. Phys. Lett.* vol. 94, 012502-012503 (2009)
18. K. Yakushiji et al., "Current-perpendicular-to-plane magnetoresistance in epitaxial Co₂MnSi/Cr.Co₂MnSi trilayers," *Appl. Phys. Lett.* vol. 88, 222504-222503 (2006)
19. Z. Diao et al., "Spin-transfer torque switching in magnetic tunnel junctions and spin-transfer torque random access memory," *J. Phys.: Condens. Matter* vol. 19 165209 (2007)
20. Y. Chu et al., "Electric-field control of local ferromagnetism using a magnetoelectric multiferroic," *Nature Materials*, vol. 7, pp. 478-482 (2008)
21. N.A. Spaldin, R. Ramesh, "Electric-field control of magnetism in complex oxide thin films," *MRS Bulletin*, vol. 33, pp. 1047-1050 (2008)
22. S. Datta, B. Das, "Electronic analog of the electro-optic modulator," *Appl. Phys. Lett.* vol. 56, pp. 665 (1990)
23. K.C. Hall, M.E. Flatte, "Performance of a spin-based insulated gate field effect transistor," *Appl. Phys. Lett.*, vol. 88, pp. 1-3 (2006)

Chapter 3.

Overview of Magnetic Coupled Spin-Torque Device (MCSTD)

3.1 Introduction

Despite the discovery of a number of fundamentally new spintronic phenomena and major progress in our understanding of the basic physics, we are still far from demonstrating useful logic devices which take advantage of spintronics. For example, for field effect transistors, which utilize spin currents, there are significant challenges in spin-injection, spin transport and spin detection and, furthermore, in devising schemes that would enable them to compete in performance with conventional charge based devices. For these reasons, several “non-transistor” based spintronic logic devices have been proposed [1-7]. We present a new spintronics device architecture using Magnetic Coupled Spin-Torque Devices (MCSTD) [8,9] that falls into the second category in device architecture, but operates by modulating the energy barrier needed to change the state of a device as in the first category. The MCSTD logic device has power gain and fan-out, and can implement the entire Boolean logic family of devices. In this chapter, we discuss the operation mechanism of the device.

3.2 MCSTD Structure and Operation

With the advent of nanotechnology, we have gained new insight and greater control over magnetic moments. Nano-fabrication processes, such as e-beam lithography and etching techniques allow us to build magnetic devices with critical dimensions well below 50 nm. This opens up a new way to control magnetizations because magnetic devices at this dimension can accommodate only a single magnetic domain. Furthermore, at these dimensions, the magnetization of the device can be electrically switched by a phenomenon called spin-torque transfer [10-14]. Spin-torque devices consist of a pair of ferromagnetic metal layers, the free and fixed layers, separated by a non-magnetic metal or a tunneling oxide layer. When electric current is passed through the device, the magnetization of the free layer is switched by spin-angular momentum transferred from spin-polarized electrons. Thus, nano-dimensioned spin torque devices, such as spin-valves and magnetic tunnel junctions (MTJ) have the potential of becoming ideal building blocks for an entirely new approach to non-volatile logic. In this section, we

describe magnetic simulations that illustrate the basic device concepts and design parameter space.

The most useful way to understand these devices is to view the energy profile of a spin-torque device. A spin-torque device with an elliptical cross section has magnetic shape anisotropy that tends to align its magnetization along the long axis. As a result, there are two stable states where the magnetization is either pointing ‘up’ or ‘down’ as shown in Fig. 3.1. These two states can be mapped to logical ‘0’ and ‘1’ for logic operations. There is an energy barrier separating the two states, which has to be overcome in order to switch the magnetization. One possible way of building a magnetic logic device is by controlling this energy barrier with input signals. Components that set this barrier are, (1) magnetic anisotropy (shape or crystalline), and, (2) external magnetic field. While both of these can be used to control the energy barrier, we focus on utilizing coupled magnetic fields from one or more adjacent devices to modulate the energy barrier of the switching device in our proposed architecture.

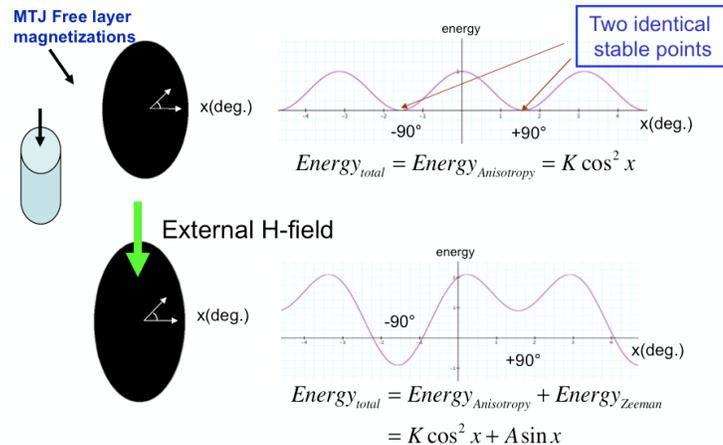


Figure 3.1. **Energy barrier of nano-pillar shaped spin-torque device.** The elliptical shape induces two stable states at +90° and -90° magnetizations.

One common way to generate magnetic fields is to induce them using current. Logic devices that switch their magnetizations with current-induced magnetic fields have been proposed in previous work [7]. However, these devices are not favorable to device scaling for the same reason as conventional Stoner-Wohlfarth MRAMs: the current and

hence power to generate the magnetic fields that can address a single bit increase with device scaling.

Our solution to generate magnetic fields that is compatible with device scaling comes from using additional spin-torque devices. If two additional devices are placed in the proximity of a few tens of nm to the original device (Fig. 3.2), their fringing magnetic fields are sufficiently strong to couple to the device and control its switching. This is a low-power solution because, the fringing fields are “free”, i.e., no power is consumed to generate them, which is different from devices dependent upon current induced magnetic fields.

We call these additional spin-torque control devices *input devices* and the center switching spin device *the output device*: the fringing fields from the input devices induce a change in the energy barrier height of the output device enabling it to predictably and reliably switch.

Figure 3.3 is a simulation showing the strength of the fringing fields when the input devices are aligned vs. anti-aligned. When the input devices have magnetizations in the same direction, their fringing fields meet in the middle and add together. The net fringing field at the center of the output device ranges between 1000~3000 A/m or 13~38 Oe, when the distance between the input devices varies between 85~125nm. Considering that the coercivity of spin-torque devices ranges from a few tens to hundreds of Oe, the fringing fields clearly affect the switching behavior of the output device. It is also possible to engineer the coercivity fields within the range where the fringing fields of the input device produce the desired output effect.

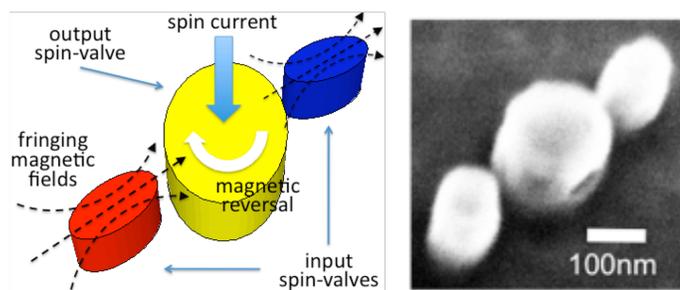


Figure 3.2. **Schematic and SEM photo of our Magnetic Coupled Spin-torque Device (MCSTD).** The MCSTD consists of two input and one output coupled spin-torque devices

When the input device magnetization changes, the net fringing fields also change. This is also confirmed in Fig. 3.3. When two input devices have opposing magnetizations, they cancel each other and the net fringing field at the center of the output device is close to zero, leaving the magnetic barriers unchanged.

Fig. 3.4 shows the energy barrier changes of the output device with fringing field

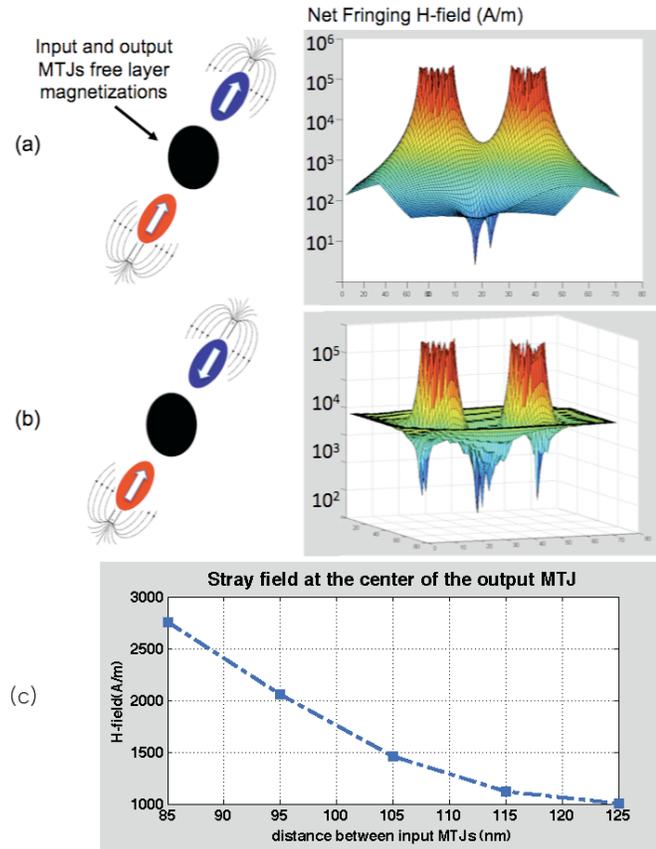


Figure 3.3. **Net magnetic fringing fields** (a) input device magnetization in the same dir. (b) opposite dir. (c) net fringing fields at the center of the output device

change from the input devices. This was confirmed by micromagnetic simulations, as shown in Fig. 3.4. Let's assume the magnetization of the output device is pointing 'down', here along the -90 or 270 degree direction. If both input devices have their magnetizations in the 'down' direction as in (0,0) input case, the barrier that blocks switching of the output device to 'up' i.e. here the +90 degree direction is raised, which prevents switching. In case (1,1) input, with the input device magnetizations pointing 'up', the barrier has been lowered. Finally, in (1,0) and (0,1), there is no change in the

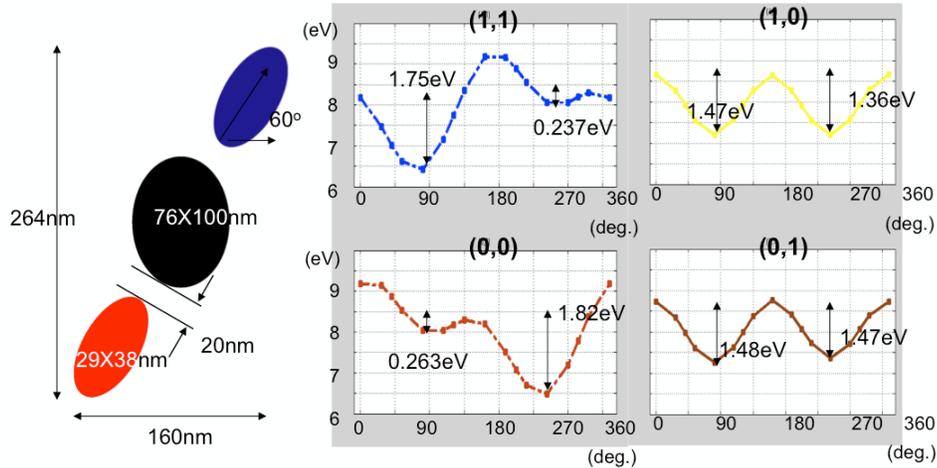


Figure 3.4. Example design of MCSTD and its energy barriers at different input MTJ magnetization configurations. ‘1’ indicates the input MTJ magnetization upwards and ‘0’ indicates downwards

energy barrier if the two input magnetizations are oriented in opposite directions. In short, depending on the input device magnetization directions, the energy barrier height of the output device can be significantly changed.

The magnetic energy of a typical spin-torque device consists of exchange, anisotropy, Zeeman and demagnetization energy, etc, which is expressed in Fig. 2.12. In making the logic device, we manipulate the magnetic anisotropy and Zeeman (external magnetic field) energy to modulate the energy barrier. Magnetic shape anisotropy and Zeeman energy components at (0,0) inputs are illustrated in Fig. 3.5. By adding the two energy components, we get the total energy in (c), an energy profile similar to Fig. 3.4(a)

This change in the energy barrier height leads to a change in the switching voltage of

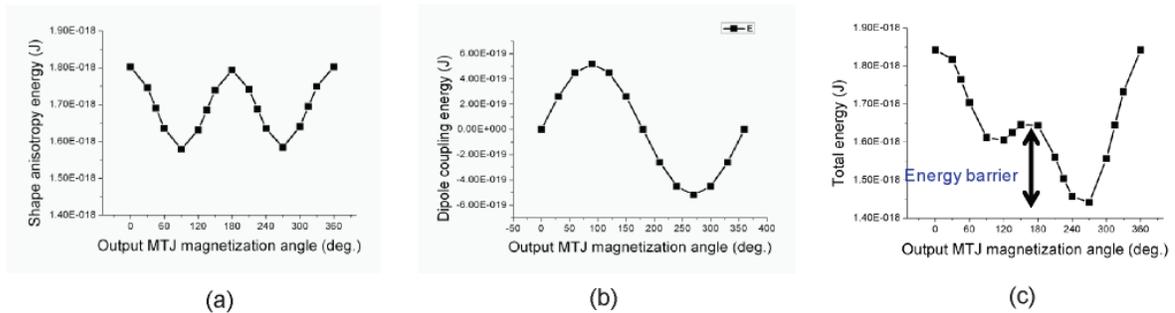


Figure 3.5. Simplified magnetic energy model of MCSTD device. For preliminary design, (a) magnetic shape anisotropy and (b) Zeeman energy (an external magnetic field at 270 deg.) are considered. Sum of (a) & (b) results in (c) total magnetic energy of the output MTJ.

the output device depending on the input device magnetization: input MTJ magnetization can be electrically controlled by spin-torque transfer effect. Input devices work as biasing dots that can make the output MTJ easier to switch or harder to switch when current is sent through the output MTJ. Energy barrier profile will change from one configuration to another as shown in Fig. 3.6 by flipping the magnetization of the input MTJs. Please note that it is the current that switches the output MTJ. Input MTJs work only as biasing fields and do not switch the output MTJ by themselves.

Using the change in the switching voltage has some resemblance to FLASH memory devices. In FLASH memory, data is stored as a change in threshold voltage, V_{th} . Electrons are injected into floating gate in order to modulate V_{th} . MCSTD is different from FLASH memory in a sense that there are two input MTJs that can individually receive signals from outside of the device. These two input devices enables two-input logic operations such as two-input NAND or NOR, which is not possible with a single FLASH memory device.

3.3 Functional Logic Design

The entire set of Boolean logic can be implemented with MCSTD gates. As shown in Figs. 3.7 and 3.10, different logic functions, i.e., NAND and XOR can be implemented by simply changing the location and the angle of the input devices. This is because different dipole couplings, hence, different energy barrier heights are produced from different locations and orientations of the input devices. This is very useful for two reasons. First, the same spin-valve (or MTJ) layer stack can be used for the fabrication of all logic functions. Different logic functions are all realized by different locations of the input devices which are realized by lithography and etching. Second, MCSTDs have a multi-dimensional design space. This is a huge improvement over other approaches to non-volatile logic, such as Magnetic Quantum Cellular Automata (MQCA) [3,4,5] where all magnetic dots have to be horizontally or vertically aligned to allow the use of antiferromagnetic coupling between the devices. There is no such restriction in MCSTD: the input devices in MCSTD can be spatially located at any point and with any angle. Magnetic coupling between devices can even be out of the plane using perpendicular

magnetic anisotropy (PMA) materials. This larger design space for MCSTDs makes the design of non-volatile logic much simpler and far more flexible.

3.3.1 NAND and NOR gates

In Fig. 3.7(a), the magnetized downward input devices ((0,0) input) raise the energy barrier to go from the ON (-90 deg.) to the OFF (+90 deg.) state. When the spin-torque current is applied to the output MTJ and the current attempts to rotate the magnetization, this rotation is blocked by the barrier, and the output device will stay in the ON state, which is the correct result for (0,0) input to NAND or NOR gates. In Fig. 3.7(b), we have (1,1) input and the energy barrier is mirror imaged to that of (0,0) input. Now the barrier is smaller and the output device can be switched to the OFF state, which is also the correct result for (1,1) input. In Fig. 3.7(c), the (1,0)/(0,1) inputs give fringing fields that cancel each other and the energy barrier is the same as the output device alone. If its barrier is large, the output device will not switch, which works as a NAND gate. If it had a small barrier, it will be able to switch, giving a NOR gate. In summary, the difference between MCSTD NAND and NOR gates is the original energy barrier height of the output device, which is determined by the spacing to the input devices or the aspect ratio of the output MTJ. Fig. 3.8 illustrates the estimated energy barrier heights for different input signals through micromagnetic simulation.

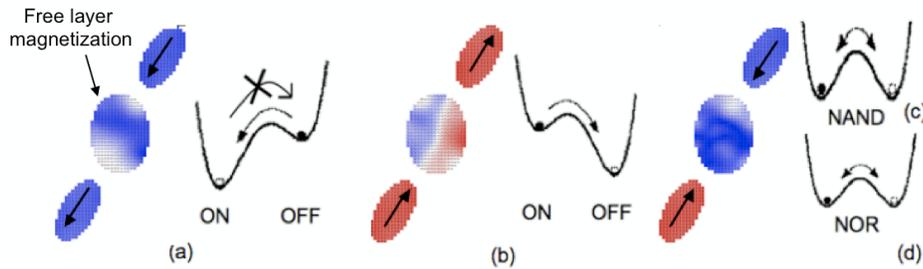


Figure 3.7. **MCSTD NAND and NOR gate concept.** Magnetization of the free layer (left) and energy barrier height profile (right) are shown. Blue color indicates magnetization downwards and red indicates magnetization upwards

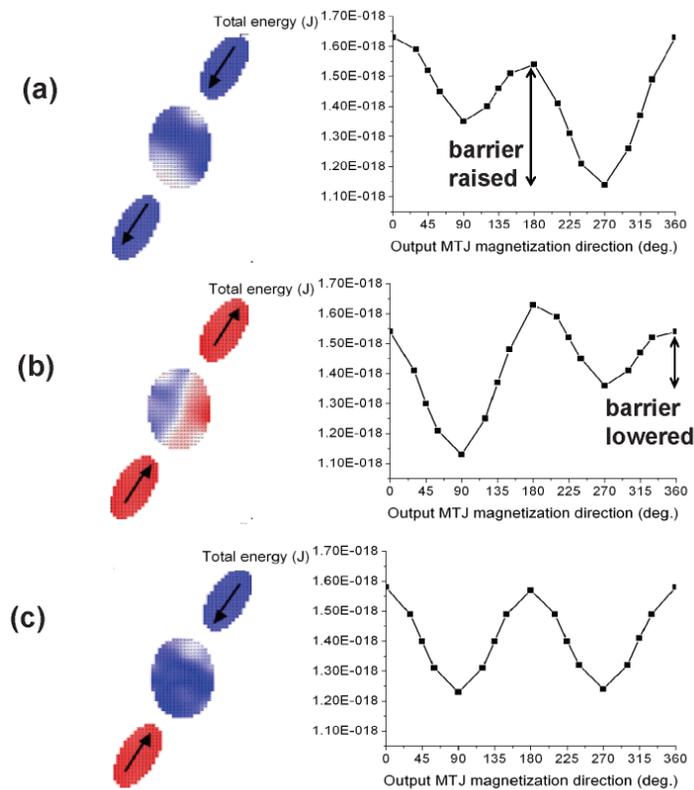


Figure 3.8. **MCSTD NAND gate energy barrier height versus input device magnetization directions.** When the initial magnetization in the output MTJ is pointing 270° , (a) (0,0) input (represented by the magnetizations in the input devices) induces the largest barrier. (b) (1,1) input, barrier lowered. (c) (0,1) and (1,0) input, no barrier lowering. All energy barrier heights are obtained from micromagnetics simulations

Figure 3.9 shows the time evolution of a MCSTD NAND gate with three different inputs. Switching of the output MTJ is gated by the input MTJ magnetizations and

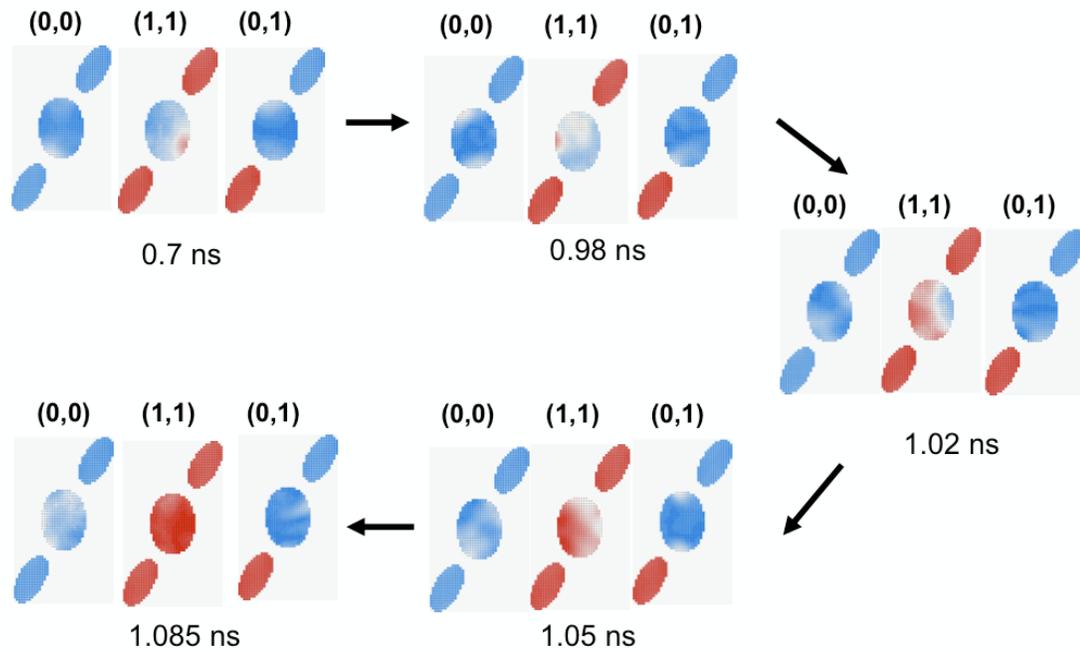


Figure 3.9. Red : spin up, Blue : spin down, $J=3.63 \times 10^7$ A/cm², switching happens at (1,1)=(spin_up, spin_up) input only

happens only for (1,1) inputs at 1.085ns ($J=3.6363 \times 10^7$ A/cm²). This happens because the energy barrier is lowered only for (1,1) inputs.

3.3.2 XOR and XNOR gates

By aligning the input devices along the hard axis of the output device as shown in Fig 10, the MCSTD gate performs a XOR operation. When the input devices are magnetized along the horizontal axis by a (1,1) input (Fig. 3.10(a)), the fringing fields will be added to lower the barrier at $x=0^\circ$ (Note that positive horizontal axis is defined as $x=0^\circ$). Now the switching of the output device is possible and the output device will switch from the ON state to the OFF state, which is the correct result for XOR operation. If the input devices are magnetized in the -x-direction by a (0,0) input as in Fig. 3.10(b), the barrier at $x=180^\circ$ will be lowered and the switching will also be possible. In contrast, if the inputs are (1,0) or (0,1) as seen in Fig. 3.10(c), the fringing fields cancel each other and the initial energy barrier of the output device has not been lowered. Therefore the switching will not happen and the output device will stay in the “ON” state, which is the correct

result for a XOR gate.

It is a significant advantage that MCSTD takes only one gate to perform the XOR operation: XOR operation is usually considered ‘expensive’ because, for example, CMOS requires 16 MOSFETs to perform this function. This compact XOR MCSTD gate simplifies full adder design (Chap. 5.1) and logic embedded biosensor (Chap 5.2) applications.

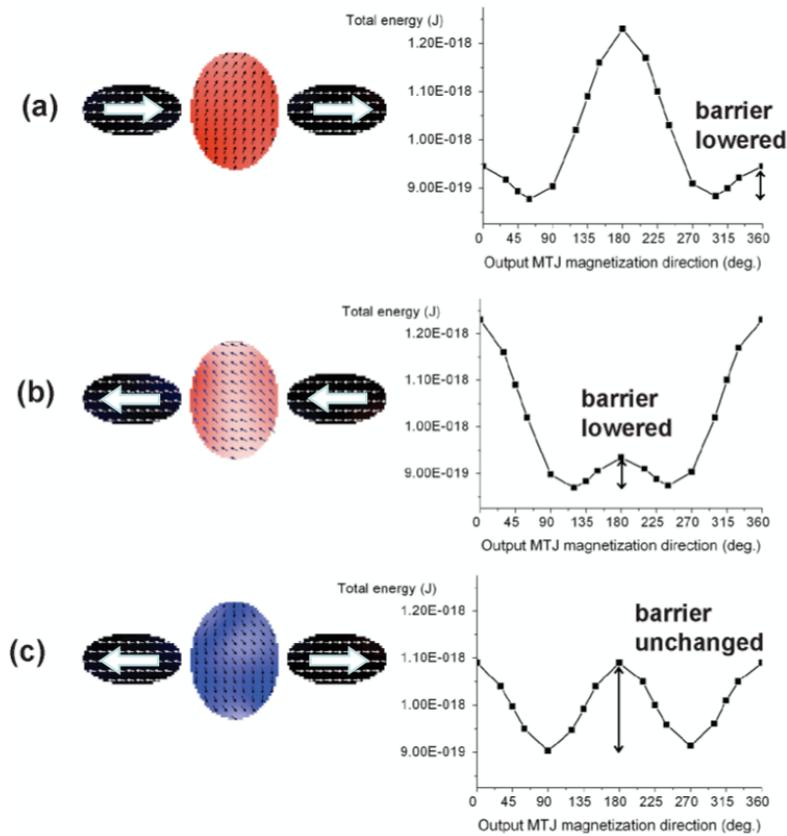


Figure 3.10. **MCSTD XOR gate energy barrier height vs. the input MTJ magnetizations.** When the initial magnetization in the output MTJ is pointing 270°, (a),(b) (0,0) & (1,1) inputs (represented by the magnetizations in the input MTJs) lower the energy barrier and facilitate the switching. (c) (0,1) and (1,0) input, no barrier lowering. All energy barrier heights are obtained from micromagnetics simulations [OOMMF]

3.4 Design considerations

3.4.1 Voltage shift

Shift in switching voltage in MCSTD gates can be explained by equation (1) and Fig. 3.11. The stray magnetic fields from the input MTJs add to or subtract from the coercivity field, H_c and shift the switching voltage points. Magnetic coupling inside a MCSTD gate makes a stronger impact to the switching voltage when the coercivity of the output MTJ is comparable to the stray fields from the input MTJs. As shown in Fig. 3.11, measured H_c of this particular output MTJ (without the input MTJs) is 75 Oe. Net fringing magnetic field from the input MTJs can be 18 Oe, which is a reasonable amount from the micromagnetics simulation result shown in Fig. 3.3(c). In this case, the fringing field amounts to 24% of H_c and results in a switching voltage shift that is much greater than thermal noise and can be utilized for logic applications.

The stray field decreases with d^3 (see equation (1), d : distance in space. Stray field decreases with d^2 , between two points in plane) as the magnetic flux density of a magnetic dipole decreases with d^3 . This imposes a fabrication challenge to close the gap between the input and the output MTJ to $\sim 20\text{nm}$. On the other hand, magnetic coercivity, H_c of the output MTJ depends on the anisotropy field, which is related to the aspect ratio

of

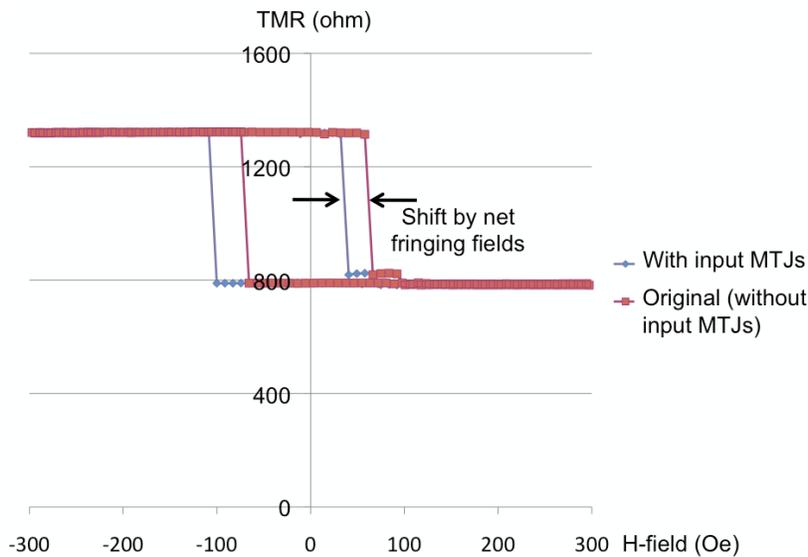


Figure 3.11. **The shift of MR vs. H-field loop** due to net fringing fields from the input MTJs. Net fringing fields work as an external magnetic fields to shift the magnetic switching points

the output MTJ : smaller aspect ratio makes H_c small. Figure 3.12 (a),(b) shows an interesting case, when magnetic coupling is increased further by making the gap 6 nm and reducing the output MTJ aspect ratio. As shown in Fig. 3.12 (a),(b), there is only one energy barrier instead of two: as magnetic coupling increases, the smaller barrier disappears and MCSTD gate has a single stable state. In this case, the output MTJ can switch without any driving current. As soon as the input MTJs switch, the stable state location changes and the output MTJ magnetization naturally relaxes to the stable point. Thus, there are two different switching mechanisms in MCSTDs depending on how many stable states it has. Hence there are bi-stable state and uni-stable state MCSTDs.

$$I_{threshold} \approx (H_k + 2\pi M_s \pm H_{stray})(4r^2t)M_s \frac{\alpha}{\eta} \frac{2e}{\hbar} \quad \dots (1)$$

$$B = \frac{\mu_o m}{4\pi d^3} (\widehat{a}_R 2 \cos \theta + \widehat{a}_\theta \sin \theta) \quad \dots (2)$$

$$\tau_{switching} = \frac{\ln(\pi / 2\theta_o)}{|I - I_{co}|} \quad \dots (3) \quad I_{wire\ current} = \frac{crH_{Oersted}}{2} \quad \dots (4)$$

$\left(\begin{array}{l} M_s \text{ saturation magnetization, } H_k \text{ anisotropy field, } r \text{ radius,} \\ t \text{ thickness, } \alpha \text{ LLG damping coefficient, } \eta \text{ spin polarization} \\ \text{factor, } m \text{ magnetic dipole moment, } \tau \text{ MTJ switching time,} \\ \theta_o \text{ the angle between the free \& fixed layers of the input MTJ} \end{array} \right)$

3.4.2 Output MTJ driver current and device reset

One of the drawbacks of using MTJs for logic application is its poor I_{on}/I_{off} ratio. If connected to a constant voltage source, the MCSTD gates will have continuous leakage current because of their small ON/OFF state resistance ratio, e.g., $I_{on}/I_{off} = 4.5$ if TMR=350%. Our solution is to provide the bias current to drive the output MTJ with a

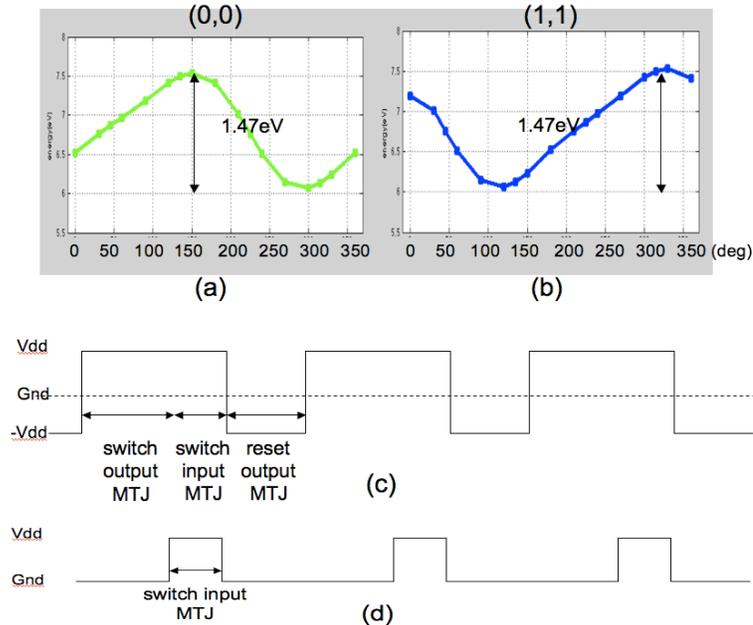


Figure 3.12. **Energy profiles of uni-stable state MCSTD and clocking schemes.** When magnetic coupling is sufficiently large, MCSTD gets single stable state ((a),(b)). The clocking schemes are different for bi-stable state MCSTD (c) and uni-stable state case (d)

clock signal. Output MTJ driver current becomes zero, when the clock signal goes down. All MCSTD gates become disconnected from voltage sources when no output MTJs need to be switched. Although no current flows, the MCSTD circuit retains the machine state due to its non-volatility. Pulse width and swing is different for two MCSTD operation mechanisms mentioned in Sec. 3.3.1. The output MTJ driver current pulse width is shorter in the single stable state case (Fig. 3.12(d)) as the output MTJ does not need to be switched by the bias current.

Bi-stable MCSTD gates have to be reset at every clock period akin to the precharge of dynamic CMOS circuits. Reset can be done by flowing the output MTJ driver current in the opposite direction from that of the switching mode. Bias voltage swings from V_{dd} to $-V_{dd}$ for the reset. A reset is not required for the uni-stable MCSTDs.

The input MTJs are tilted to give some initial angle between its free and the fixed layer to accelerate its switching as shown in equation (3).

3.4.3 Oersted field induced noise

Magnetic reversal is an interaction between the spin torque momentum and the Oersted field. The Oersted field is a concentric magnetic field generated when current is applied in a normal direction to the surface. On the other hand, spin-torque momentum is unidirectional. Due to the difference in the field shape of the two, the Oersted field can

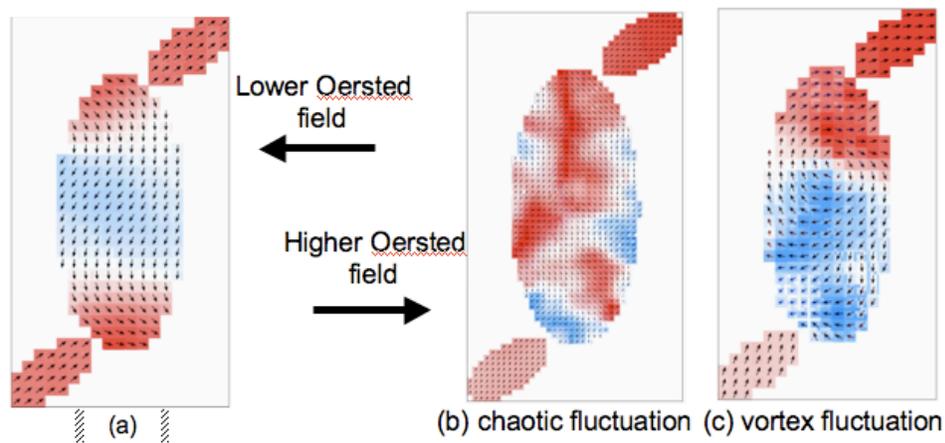


Figure 3.13. Influence of the Oersted field on spin-torque switching process

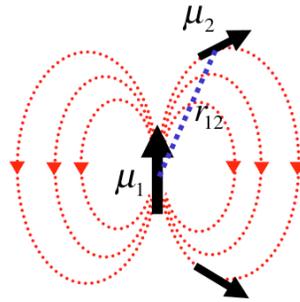
interfere with the spin torque momentum. Actually, the Oersted field and its interference with spin-torque facilitate the magnetic reversal in the early stage of switching. But, constant interference with spin torque generally delays the reversal process. When it happens random and chaotic fluctuations seen in Fig. 3.13 (b), (c) occur and magnetic reversal gets delayed. Equations (1) and (4) show that the Oersted field decreases linearly with device radius r when current density is maintained constant. This is one reason why device scaling is favorable for MCSTD: the other reason is power consumption.

3.4.4 Positions and angles of the input spin-torque devices

Different positions and angles of the input spin-torque devices were compared to find the largest dipole coupling. Dipole coupling represents the strength of the fringing fields, whose general form is shown in Fig. 3.14.

To begin with, dipole couplings of two magnetic moments are compared (Fig. 3.15). When the magnetic moments are aligned along their longer axis, the energy is lowest when the magnetic moments are facing the same direction than when they are opposing each other. On the other hand, when the magnets are beside each other, the system become more stable, i.e., lower energy when they are anti-aligned. Next, we extend the discussion to three nanomagnets. The lowest energy state is when each pair of nanomagnets are in the lowest energy states individually. This leads us to the conclusion that the dipole coupling energy is the lowest when the three magnets are all aligned along the longer axis. For complete analysis the input MTJ location and angles are varied as shown in Fig. 3.17 to get dipole coupling strength for various conditions. In the case when the input and output MTJs are aligned along the longer axis (90° case in Fig. 3.18), dipole is maximum but it can lead to lower areal density because the distance between the neighboring MCSTD gate gets shorter. Also, the layout in crossbar architecture becomes difficult. Due to these reasons, the input MTJs are therefore aligned to be slightly off the longer axis line and at angles of $60\sim 80^\circ$.

As was discussed in previous sections, the MCSTD gate is based on its “input signal dependent switching voltage” for logic operations: the switching voltage of a MCSTD gate becomes different depending on the input signal. Figure 3.19 shows the proposed switching probability versus voltage characteristics of NAND and NOR gates. As we



$$\text{dipole coupling} = \sum_{i,j \neq i} \frac{\vec{\mu}_i \cdot \vec{\mu}_j}{r_{ij}^3} - \frac{3(\vec{\mu}_i \cdot \vec{r}_{ij})(\vec{\mu}_j \cdot \vec{r}_{ij})}{r_{ij}^5}, \mu = \int M dV$$

Figure 3.14. Dipole coupling between the input and output MTJs. μ represents the magnetic moment and r is the distance between the magnetic moments.

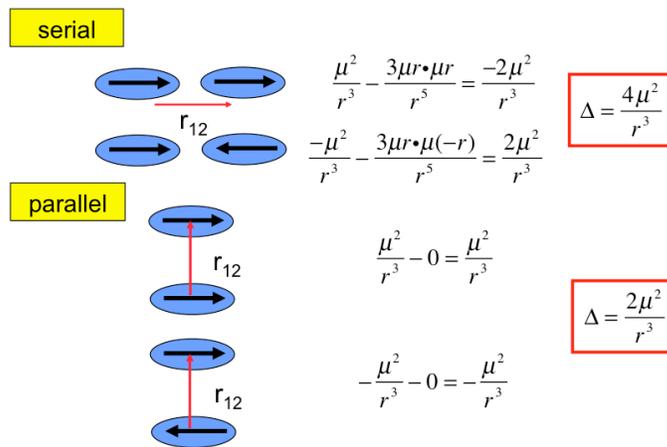


Figure 3.15. Dipole coupling between two magnetic moments

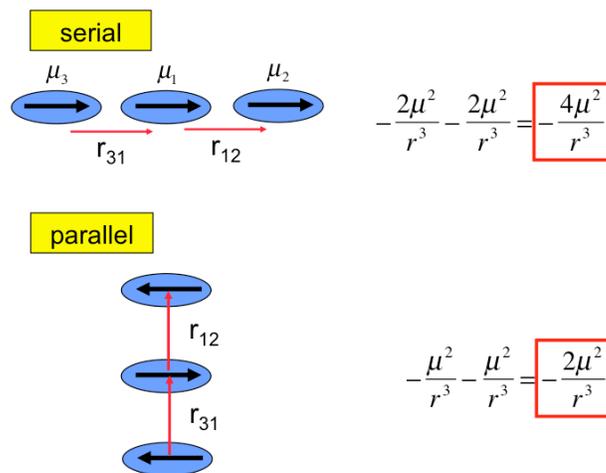


Figure 3.16. Dipole coupling among three magnetic moments

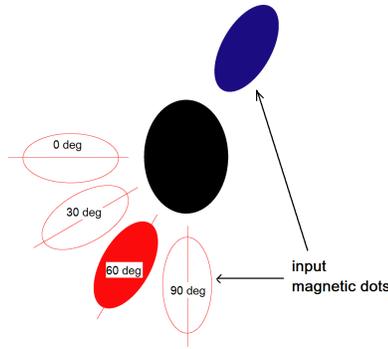


Figure 3.17. MCSTD gate with the input MTJ at different locations and angles

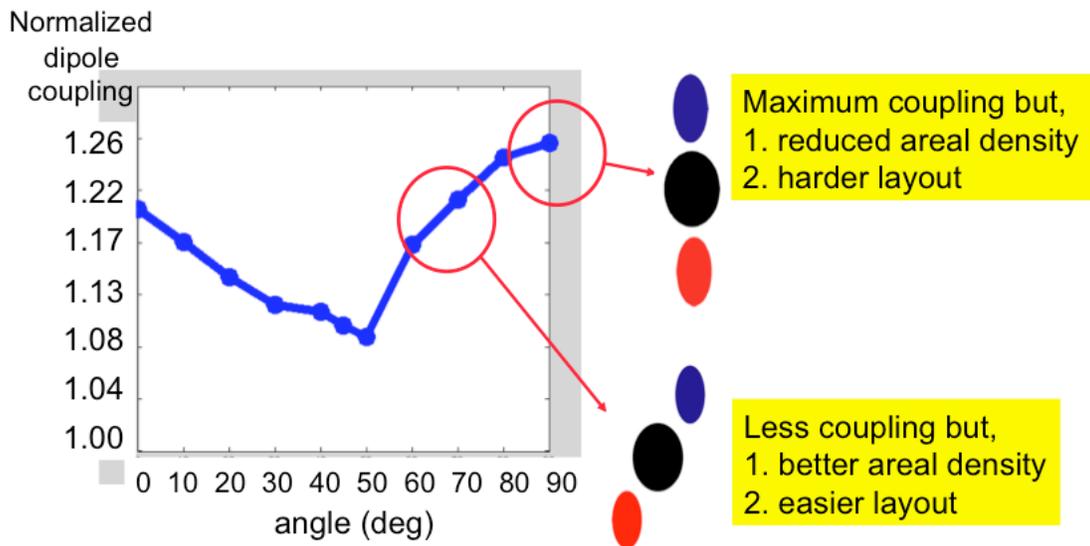


Figure 3.18. Dipole coupling strength vs. input MTJ angles

sweep the biasing voltage, the MCSTD gate with a low energy barrier height switches at a lower voltage and the other ones with higher barriers switch at larger voltages. When a MCSTD gate switches, it goes from the ON state to the OFF state. For NAND logic, we want the MCSTD gate to switch only at (1,1) input, which requires the switching probability plots for (1,0) and (0,1) inputs to shift closer to or come after the (0,0) input. In contrast, for NOR gate, (1,0) and (0,1) inputs should switch at a smaller voltage or a closer voltage to that of (1,1) input.

Changing the location and angle of the input MTJ controls the net fringing field direction and hence the switching voltages at (0,1) and (1,0) inputs. Let's take a look at the energy barrier heights of (0,0) (Fig. 3.20 (a)) and (0,1) (Fig. 3.20 (b)) input cases for

the comparison of their switching voltages. Usually, the energy barrier heights for (1,0) and (0,1) inputs come between that of (1,1) and (0,0) inputs. This is because the fringing fields from the input MTJs are cancelled at (1,0) or (0,1) inputs to make virtually no change in the energy barrier height from that of the output MTJ alone. For (0,0) input, the energy barrier gets raised from that of the output MTJ. In this way, we can make a NOR gate with the switching probability plots shown in Fig. 3.19. In order to make a NAND gate, the barrier height of the (0,0) input should be comparable to or smaller than that of (1,0) and (0,1) inputs. This can be done by increasing the horizontal components of the fringing fields. As shown in Fig. 3.20 (c), the horizontal fringing fields are aligned with the magnetic hard axis of the output MTJ. Fringing fields pointing left (Fig. 3.20 (c)) will lower the barrier in 180° and increase that in 0° in the energy profile plot. From Fig. 3.20 (a) and (b), the energy barrier heights at (1,1), (0,0) and (1,0) inputs are $1.9 \times 10^{-19} \text{J}$, $4.0 \times 10^{-19} \text{J}$ and $3.3 \times 10^{-19} \text{J}$ respectively. The device shown in Fig. 3.20 (c) performs NAND operation because the energy barrier heights for (1,0) and (0,1) inputs are closer to that of (0,0) inputs than (1,1) input. The ways to decrease the horizontal component (and to make NOR gate) is to 1) align the input MTJs along the output MTJ easy axis 2) tilt the input MTJ angle as it is shown in Fig. 3.21 (b). Fig. 3.21 shows the fringing magnetic fields in the device. The colored regions (red and blue) represent the horizontal components. The horizontal components in the fringing fields have decreased with the above methods.

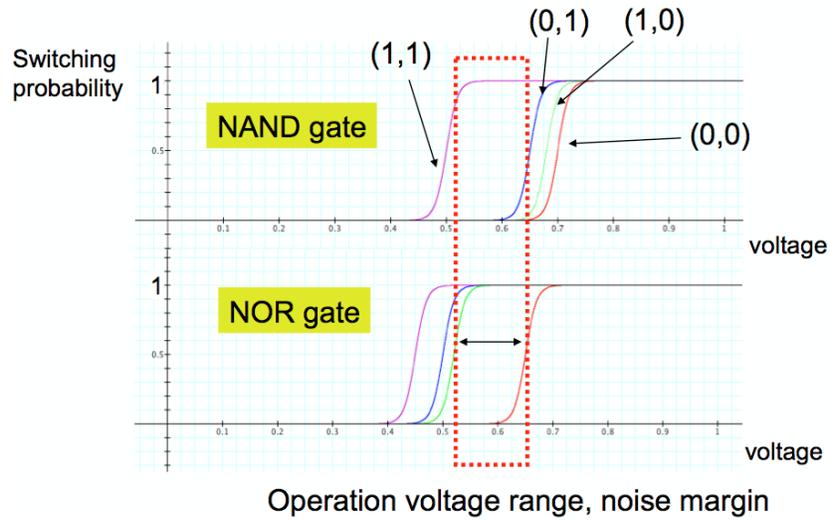


Figure 3.19. **Proposed voltage operation range of MCSTD gates.** NAND and NOR gates have different switching voltage for different inputs, which makes them function as different logic devices. The numbers in parentheses are input signals.

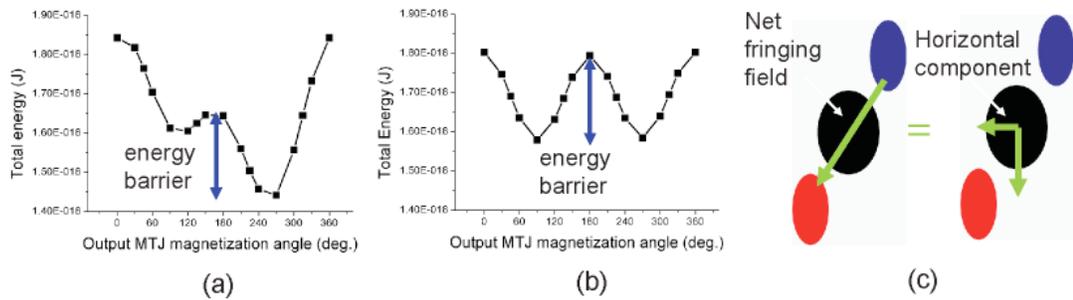


Figure 3.20. **MCSTD NAND gate energy barrier heights.** Energy barrier at (a) (0,0) input (b) (1,0) input

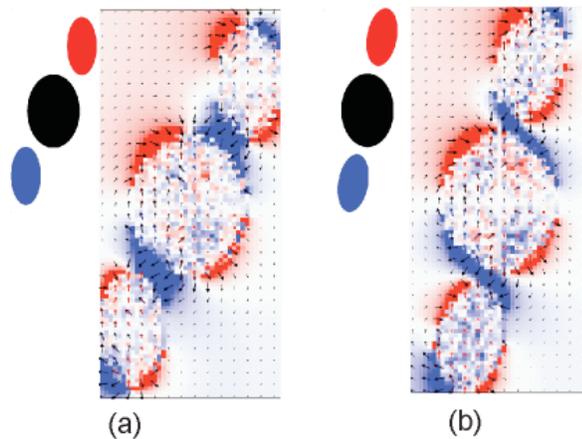


Figure 3.21. **Horizontal components of fringing magnetic fields** (Blue: left, Red: right) (a) NAND (b) NOR gate

3.4.5 Materials for the input device

Now the material for the input device was varied as follows in order to see their impact on the switching characteristics. NiFe is used for the output MTJ free layer for all cases.

a. NiFe

The same material was used for both the free layer of the input and the output MTJs. A MCSTD shows different switching speed for the higher and lower barrier case set by the input device magnetizations. The speed difference becomes more than three times as the current density to drive the output spin-torque device gets smaller than $2 \times 10^7 \text{ A/cm}^2$.

b. Co

Co is a “harder” magnetic material than NiFe. The difference in switching speed becomes much greater. At $J=1.8 \times 10^7 \text{ A/cm}^2$, the lower barrier case switches within 2ns, but the higher barrier takes more than 10ns to switch. When $J=1.6 \times 10^7 \text{ A/cm}^2$, the higher barrier case does not switch no matter how much time is given. A larger difference in switching speed is favorable for the MCSTD to be used as a logic device as it reliably switches on and off.

c. CoFe

When an even harder magnetic material (CoFe) is used for the input device, the switching time of the lower barrier case gets even shorter. To the contrary, the higher barrier cases did not switch at all within $J < 5 \times 10^7 \text{ A/cm}^2$. This makes the switching speed difference very large for a MCSTD.

In summary, using heterogeneous materials for the input device makes the switching speed difference large and is very favorable for the MCSTD logic device. Heterogeneous material, preferably harder magnetic material for the input spin-torque device will achieve larger switching speed difference.

3.5 MCSTD Logic Design

MCSTD is a general logic device that has power gain, fan-out and signal level restoration capability that allows it to be cascaded indefinitely. These characteristics originate from the interplay among asymmetry in the input and output device size and

assistance from the external current source. In this section, we discuss these attributes in detail.

3.5.1 Gain

An MCSTD gate takes the input device magnetization as an input and generates the output device magnetization. Given that the magnetic properties and the thickness of the input and output devices are the same, the gain of the MCSTD gate becomes the ratio of between the input and the output device areas. To make the gain > 1 , the input devices are smaller in size and magnetization than the output device. In order to achieve switching the magnetically “stronger” output device with the “weaker” input device, the output MTJ driver current is supplied from an external voltage/current source. In other words, all the output MTJs are connected to an external voltage/current source, which is implemented with CMOS circuit. This is similar to a CMOS circuit, where an input signal is used to charge the gate capacitance to turn on MOSFET devices only. CMOS circuit also uses an actual current to raise the output potential from an external voltage source or supply voltage, V_{dd} .

In addition, feedback between the input and the output devices, i.e., the output disturbing its own input, is prevented by giving higher aspect ratio shape to the input devices. This increases magnetic anisotropy and stability, which avoids dipole coupling from the output MTJ switching the input MTJ free layer.

3.5.2 Nonlinearity

For a logic device to have sufficient noise margin, a nonlinearity between the output and the input signals is required. In other words, any signal that comes within the range of the noise margin of the allowed state variable signal levels, should be amplified to reach stable values. MCSTD has excellent non-linearity and its input and output signals can be modeled as a Sigmoid function shown in Fig. 3.19. If the input signal lowers the output MTJ energy barrier below the threshold within its noise margin, the output MTJ switches with the output MTJ driver current. If the barrier height stays above the threshold, the output MTJ does not switch.

3.5.3 Cascadability

Signal propagation and energy supply paths are separate in MCSTD circuits in order to acquire signal restoration capability. As shown in Fig. 3.22 (b), signals from previous stage MCSTD gates propagate through interconnects and switch the input devices. Then, the input device magnetizations bias the output device of the subsequent stage, which allows the output devices to evaluate the input signals based on their logic functionalities. However, the actual energy to switch the output device comes from the external voltage supply discussed in Sec. 3.4.1 (Fig. 3.22 (c)). If the energy barrier has been lowered by the input device magnetization, the spin-torque from the external voltage supply can switch the magnetization to the opposite level. In other words, the signal has been restored to stable signal levels at every output device without any signal degradation. This allows the MCSTD gate to be infinitely cascaded. When MCSTD use both the input and output device have the magnetization as the state variables, no signal conversion is required for cascading the devices.

When MCSTDs are cascaded, the output MTJ of one stage is serially connected to the input MTJs of the subsequent stage (as in Fig. 3.22). This forces the input MTJs of the next stage to be oriented as the spin polarity signal of the previous stage output MTJ. In addition to the serial signal passing, a bias current, actually supplying the energy that drives the next stage input MTJ is supplied from an external supply voltage so that the voltage applied to the output MTJ always has the same value regardless how many gates are cascaded. This allows the signal to propagate without degradation, i.e. a form of *signal level restoration* function. The output and the input MTJ resistance add up in series but the total resistance is within a certain range as the input MTJs form a parallel circuit with other input MTJs.

3.5.4 Fanout

A MCSTD gate can have a fanout > 1 , i.e., it can drive multiple other MCSTD gates due to the asymmetric area ratio between the output and the input devices (Fig. 3.22 (a)). The output MTJs are designed to be larger than the input MTJs, which means that the current that flows through the output device is many times larger than that required to switch the input devices in the subsequent stage. Thus, a single output MTJ can drive

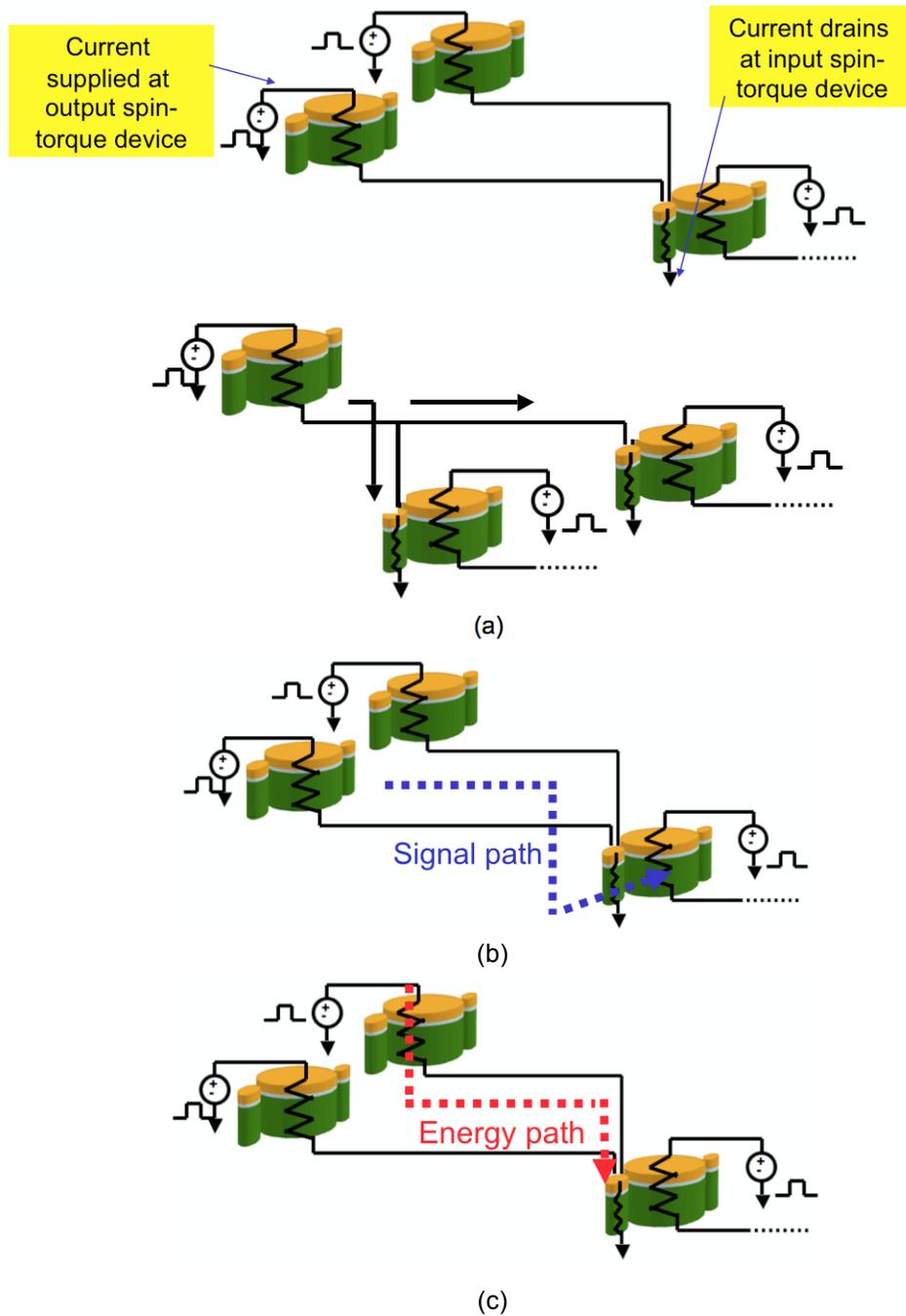


Figure 3.22. **Circuit configuration for MCSTDs.** To cascade MCSTDs, the input MTJs of the next stage are connected to the output MTJs of the previous stage. Each output MTJ is connected to an external voltage source.

multiple input MTJs of subsequent logic stages.

3.5.5 Static power consumption

Unlike CMOS circuits, MCSTD gates have “instant on/off” capability because they can retain the device states even when they are disconnected from the power supply. This is because the spin-torque devices are non-volatile as the magnetic materials in these devices are the same material (Co, Fe, Ni) as that in hard disks or magnetic random access memories (MRAMs). However the spin-torque devices usually do not have a good I_{on}/I_{off} ratio- typically less than 6x – which can lead to a huge amount of leakage current. This problem can be avoided if pulse-mode current signals are supplied to the output MTJs to drive them. The output MTJ driver current can be completely turned off when the circuit is not in use because the devices are non-volatile. Furthermore, there is no overhead in power on/off the MCSTD circuit because power-gating can be done at the chip-level. This eliminates the need for power-gating transistors to be distributed across the chip, which largely accounts for the performance reduction in power-gating techniques in CMOS circuits.

References

1. D. A. Allwood, et al., “Magnetic domain-wall logic,” *Science* 309, 1688–1692 (2002)
2. B. Behin-Aein, D. Datta, S. Salahuddin, S. Datta, “Proposal for an all-spin logic device with built-in memory,” *Nature Nanotechnology*, DOI: 10.1038 (2010)
3. D. Carlton, N.C. Emley, E. Tuchfeld, J. Bokor, “Simulation of nanomagnet based logic architecture,” *Nano Lett.* **8**, 4173–4178 (2008)
4. R. P. Cowburn, M.E. Welland, “Room temperature magnetic quantum cellular automata,” *Science* **287**, 1466–1468 (2000)
5. A. Imre et al., “Majority logic gate for magnetic quantum-dot cellular automata,” *Science* **311**, 205–208 (2006)
6. A. Khitun et al., “Spin wave logic circuit on silicon platform,” *Fifth International Conference on Information Technology: New Generations*, 1107–1110 (2008)
7. A. Ney et al., “Programmable computing with a single magnetoresistive element,” *Nature* **425**, 485–487 (2003)
8. L. Leem, J.S. Harris, “Magnetic coupled spin-torque devices and magnetic ring oscillator,” *Proc. IEDM* DOI: 10.1109/IEDM.2008.4796640 (2008)
9. L. Leem, J.S. Harris, “Magnetic coupled spin-torque devices for nonvolatile logic applications,” *J. Appl. Phys.* **105**, 07D102 (2009)
10. E. Chen, “Current Status and Future Outlook of STT-RAM Technology,” *18th Information Storage Industry Consortium Annual Meeting* (2008)
11. L. Berger, “Emission of spin waves by a magnetic multilayer traversed by a current,” *Phys. Rev. B* **54**, 9353–9358 (1996)
12. J.C. Slonczewski, “Current-driven excitation of magnetic multilayers,” *J. Magn. Magn. Mater.* **159**, L1–L7 (1996)
13. J.Z. Sun, “Spin–current interaction with a monodomain magnetic body: a model study,” *Phys. Rev. B* **62**, 570–578 (2000)
14. M. Tsoi et al., “Excitation of a magnetic multilayer by an electric current,” *Phys. Rev. Lett.* **80**, 4281–4284 (1998)

Chapter 4.

Spin interconnection

4.1 Introduction

Interconnections are critical elements for all low power circuit implementations of any technology. Spin-based interconnections can be far more power-efficient than charge based technologies because there is no continuous flow of electrons or charging of interconnects. In this chapter, we introduce a number of novel interconnection schemes for Spintronic logic devices. Technical challenges that we would like to address are

1. How to overcome short spin decoherence length
2. How to communicate without electrical current, i.e., transport of electron charges

As a solution, we provide 1) Complementary MCSTD based spin-interconnection scheme (Sec. 4.2) and 2) Magnetic Domain-wall based interconnection (Sec. 4.4). Discussion of each scheme is followed by specific example applications, such as a Magnetic Ring Oscillator (Sec. 4.3). While each scheme is complete on its own, two interconnection schemes can be used together in a complementary way for scalability. In other words, use Magnetic Domain-wall based scheme for short-distance communications and Complementary MCSTD based scheme for long-range communications.

4.2 Complementary MCSTD based spin-interconnection

As spin diffusion length is very short (5nm ~1000nm [1]) and spin information is easily lost at various material interfaces, it is very difficult to transport spin information from one device to another. While there exist material research efforts to increase the spin coherence time (T_1 and T_2), we introduce a circuit and architectural approach for the first time of its kind. Our solution to this spin-decoherence problem is called *complementary MCSTD* (Fig. 4.1): a pair of MCSTD gates having the same input signals, one MCSTD gate is configured to be in a high resistance state and the other in a low resistance state. This dual gate setup is used to convert a spin (up, down) signal into current amplitude (high, low) difference information. For example, if MTJs have TMR ratio of 350%, the current ratio between two gates will be 4.5:1. The output currents from each MCSTD gate propagate in two separate wire interconnects to be spin-polarized

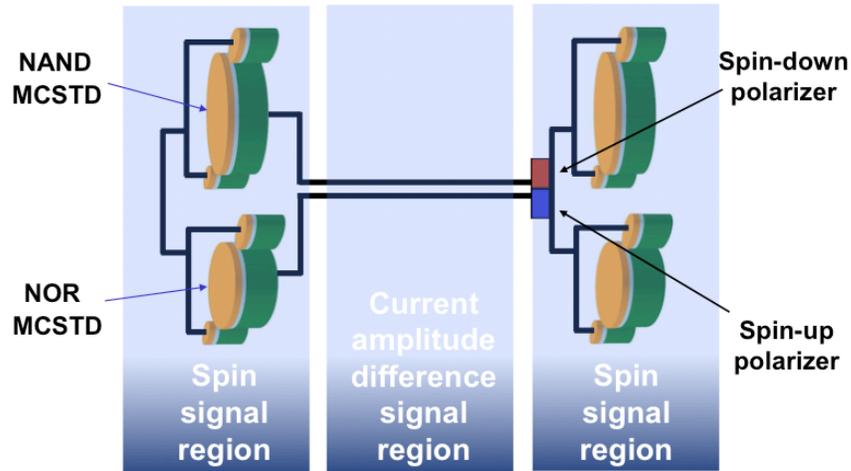


Figure 4.1. **Complementary MCSTD and spin interconnection.** Spin signal of output MTJ is converted to the current amplitude difference and reconstructed to spin signal at the next logic stage. Red and blue rectangles on top of input MTJs represent spin-up/down polarization layers. The number of wires and spin polarizer are reduced for simplicity

again and merged at the input of the next stage (Fig. 4.3). Spin-polarization information is propagated individually in two separate wires such that the current in the top wire is always spin-up polarized and the current in the bottom wire always spin-down polarized only. Since the merged current consists of spin-up and spin-down electrons in proportion to the current amplitudes in the respective interconnects, spin information is reconstructed at the input of the next stage logic devices (Fig. 4.3 (d)). For example, if each current is polarized to 90% and TMR=350%, spin polarization of reconstructed spin signal becomes 57%, which is high enough to switch the input devices in the next stage. Overall, the spin signal or the magneto-resistance information of the complementary MCSTD pair is converted into current amplitude data for communication over on-chip interconnect and then transformed back to spin information at the next stage MCSTD gates.

To make complementary MCSTDs, one can follow De Morgan's law in Fig. 4.2 akin to the dual pull-up and pull-down networks in CMOS. First, a pair of MCSTD gates implement the dual logic functions (Fig. 4.3 (a)). Second, the incoming spin polarity is inverted for one of the MCSTD gates (Fig. 4.3 (b)): spin-polarizer layout is reversed for the bottom MCSTD gate, it receives the opposite spin signals compared to the top

MCSTD gate. As mentioned in Chap. 3, all the output MTJs are connected to supply voltage, which drives them.

With this interconnection scheme, spin information can be transferred beyond the

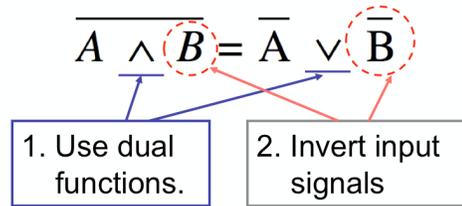


Figure 4.2. **De Morgan's law.** To implement complementary MCSTD gate, dual functions and input signal inversion are needed

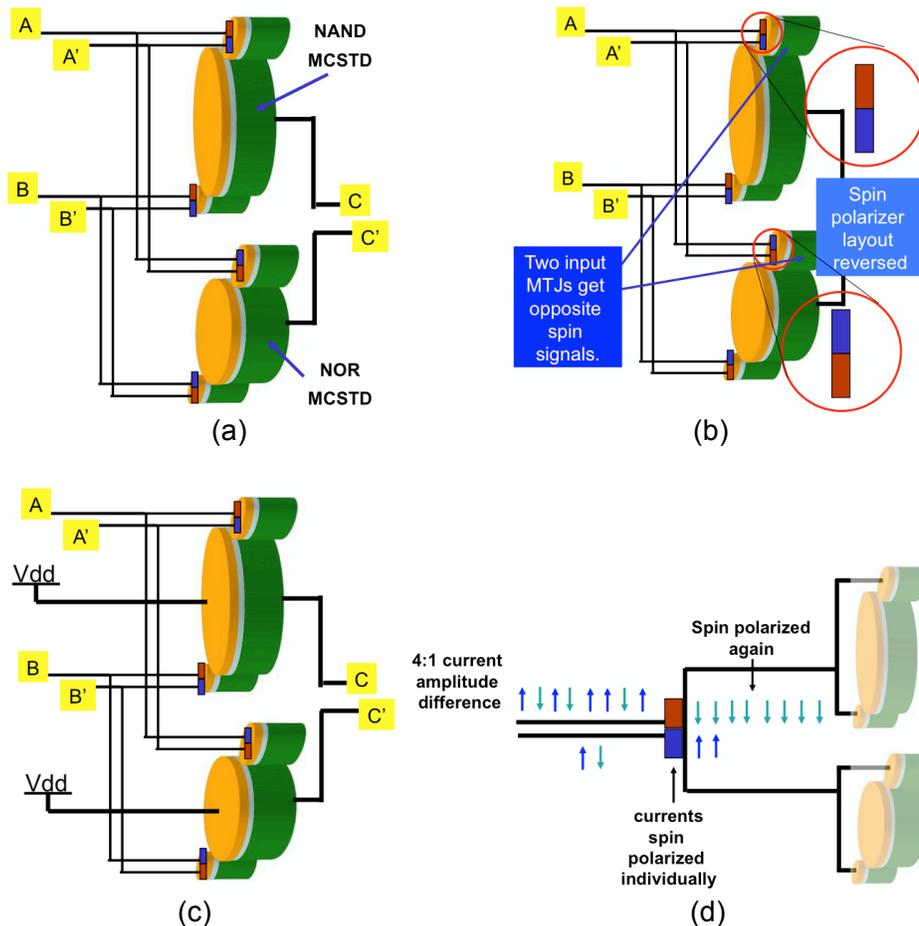


Figure 4.3. **Actual implementation of complementary MCSTD based spin interconnection.** (a) Differential signal with complementary MCSTDs (b) spin-polarizers (c) External supply voltage (d) spin signal reconstruction process. Red and blue rectangles on top of input MTJs represent spin-up/down polarization layers

short physical spin coherence length limit. MCSTDs are easily cascaded, which makes them very well suited for logic operations.

4.3 Magnetic Ring Oscillator

When MCSTDs are cascaded, the output MTJ of one stage is serially connected to the input MTJs of the subsequent stage (as in Fig. 4.4). This forces the input MTJs of the next stage to be oriented as the spin polarity signal of the previous stage output MTJ. In turn, the output MTJs are connected to external supply voltage or driver circuit, which supplies a current that flows from the output MTJ and to the input MTJ to be drained there.

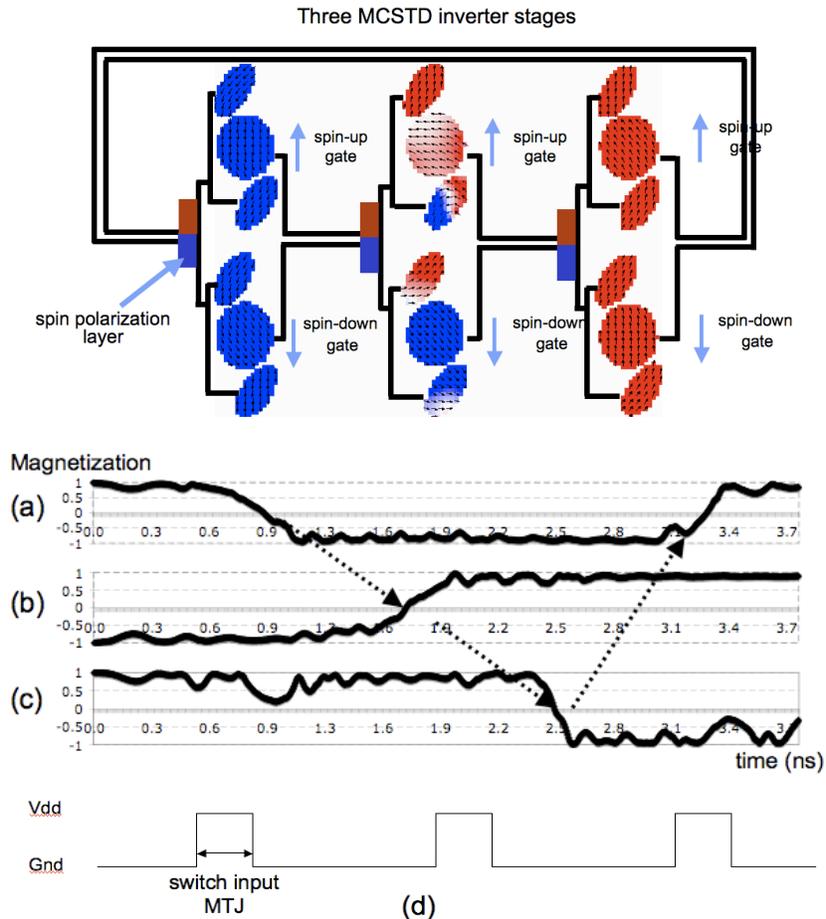


Figure 4.4. **A three stage MCSTD ring oscillator.** Fig. (a),(b) & (c) illustrate the free layer magnetization waveforms of the 1st, 2nd & 3rd inverter stage respectively. '1' spin up, '-1' spin down. Optimistic device properties of RA product=10ohm*um² and TMR=400% used for the simulation.

In order to measure the MCSTD gate delay in actual circuit configurations, a three-stage ring oscillator was implemented with the micromagnetics simulation. Each inverter stage uses a pair of uni-stable MCSTD. Two inputs of a NAND MCSTD gate are simply tied together, i.e. two input MTJs receive the same spin currents to form an inverter gate. A pair of complementary MCSTDs is used for spin-interconnection. The output MTJ driver current waveform of Fig. 4.4 (d) was useful in reducing the circuit instability. Depending on the spin interconnection design, the output MTJ can have some interference in switching direction between the one forced by the fixed layer magnetization and that induced by the input MTJs. Pulsed bias current can delay this interference until the output MTJ switches to the stable point. For these reasons, our implementation of a ring oscillator is different from simple conventional ones. Our MCSTD ring oscillator, as

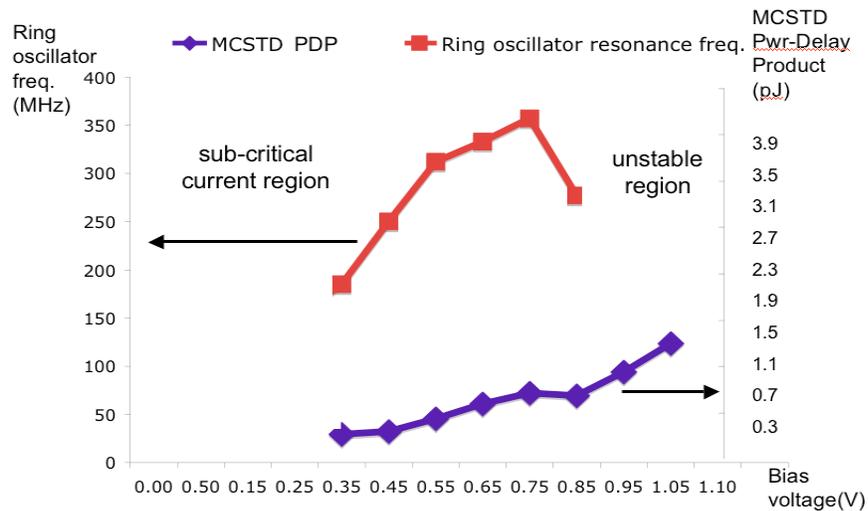


Figure 4.5. **MCSTD ring oscillator frequency and MCSTD gate energy consumption.** MCSTD ring oscillator operates within certain voltage range where current density is above the switching threshold and the Oersted field is small enough not to induce instability.

demonstrated with micromagnetic simulations (Fig. 4.5), shows a resonance frequency of 100~400MHz depending on the supply voltage. At 400MHz, each MCSTD gate shows a ~0.8ns average gate delay. It was found that there is a certain voltage range where the ring oscillator works. When the supply voltage is too low, the current density is below the critical switching current and MCSTD circuit doesn't start. As the voltage increases gates switch faster and resonant frequency improves. But at very high voltage, instability

caused by the Oersted field slows down the circuit and eventually make it fail (Fig. 4.5). MCSTDs are non-volatile devices that can work both as logic and as storage elements. With MCSTDs, no additional flip-flops or buffers are necessary to store temporary data (even under power glitches). Savings in the number of gates, energy and latency due to transferring data back and forth from temporary storage, makes MCSTDs more power efficient than conventional CMOS. Since MCSTDs are equivalent to logic embedded flip-flop, we believe implementing data-retention flip-flop with MCSTD for non-volatile logic would be a good way for the MCSTD technology to be introduced into the CMOS based semiconductor industry. The power-delay product of MCSTD is estimated to be around 300fJ (Fig. 4.5), which is 3~13x larger than that of CMOS flip-flops at the 0.20 μ m node [2]. MRAM research groups are extensively investigating ways to lower the threshold spin current density of spin-torque transfer MRAM. When the threshold spin current density decreases, the power consumption of MCSTD will be lowered by the square of its linear reduction.

4.4 Magnetic domain-wall based interconnections

4.4.1 Overview

Interconnections are critical elements for all low power circuit implementations of any technology. Non-volatile logic can be far more power-efficient than charge based technologies because there is no continuous flow of electrons or charging of interconnects. We propose using magnetic domain-wall motion for short-range (<500nm)

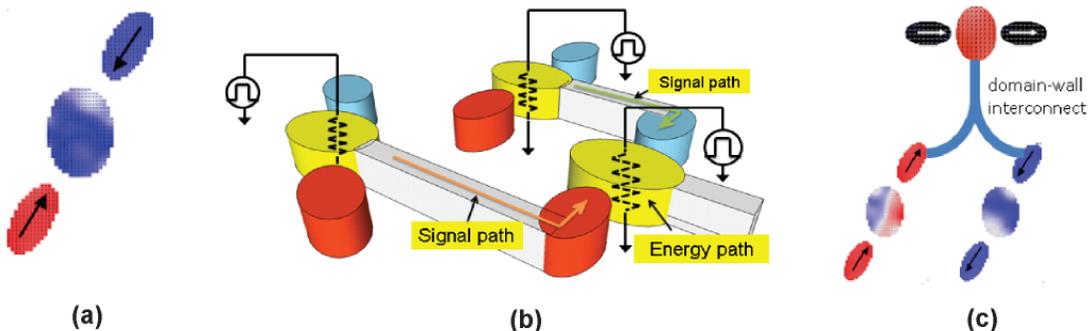


Figure 4.6. **Magnetic domain-wall based interconnections.** (a) signal gain, magnetically “weaker” input device switches the “stronger” output device (b) The output and input MTJs connected with magnetic nanowire. Signal path and energy path in MCSTD circuit for signal level restoration (c) fanout, magnetically “stronger” output device can drive multiple input devices

signal propagation. The basic idea of our domain-wall interconnect scheme is shown in Figs. 4.6 and 4.7. The output device and the input device to the following stage gate are connected via a magnetic nano-wire. Due to the fringing fields from the output devices, magnetic domain-walls in the nano-wire move in the direction of the output device magnetization. Since magnetic fringing fields are always emanating from magnetic devices, no additional power is consumed to drive the signal propagation. 30nm/ns/Oe [3] or higher mobility has been reported for field driven magnetic domain movements. As discussed in Chap. 3, fringing magnetic fields from the input/output MTJs range between 10~40 Oe, which can easily transport a domain wall over 300~1200nm in 1ns.

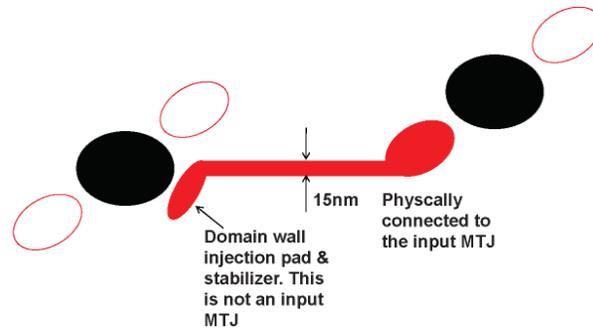


Figure 4.7. **Domain-wall motion interconnection design**

4.4.2 Reliable signal interconnection scheme

In our approach, domain-wall motion communications are made more reliable by minimizing interfering fringing fields from neighboring spin-torque devices. First, MCSTD gates are equipped with an additional soft layer on top of the free layer (Fig. 4.8) that allows us to actively turn on/off the output device fringing fields. As shown in Fig. 4.8 (a), in the normal signal-propagation mode, both the soft and the free layers are aligned in the same magnetic orientation, where two magnetic fluxes are combined to create a larger magnetic driving force for the domain-wall motion. In contrast, when the soft-layer switches to the anti-aligned mode with the free layer due to spin-torque, the magnetic flux forms a closed-loop inside the soft free layer region. Then, fringing fields from the output device will be zero and there will be no interference between the fringing fields from multiple sources. In addition, the switching current density for the output device will be reduced because whenever current is sent through the device, the soft-layer

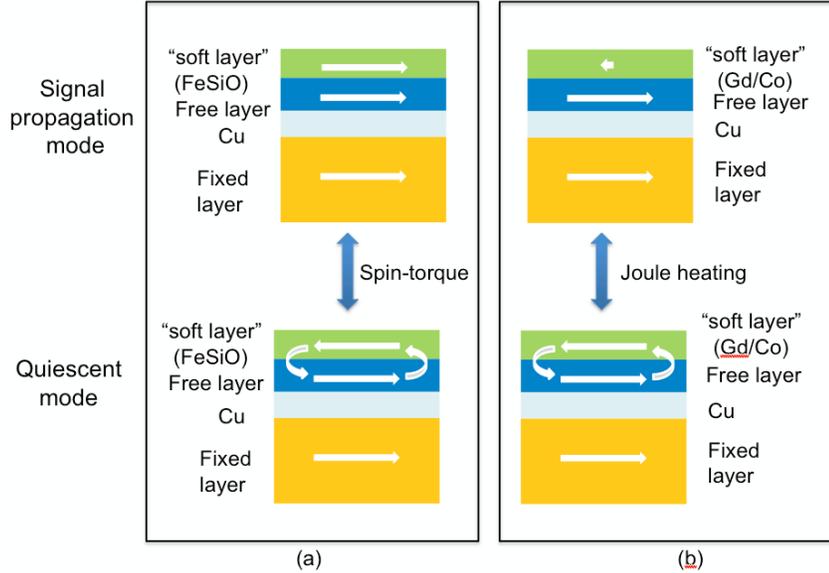


Figure 4.8. **Composite free layer structures for “interference-free” domain-wall motion interconnects.**

will switch first and facilitate the switching of the actual free layer with the exchange [5]. Another promising way to implement the normal/quiescent modes in the output device is to use a Gd/Co soft-layer that is anti-ferromagnetically coupled with the free layer (Fig. 4.8 (b)).

Second, signal propagation using domain-walls is performed in two-steps (Fig. 4.9). This is because the output devices near the both ends of the interconnect have conflicting roles and cannot be in the same mode simultaneously: while one output device that drives the domain wall is in signal propagation mode, the other should be in quiescent mode or vice versa. As shown in Fig. 4.9, when the output devices in the group 1 are propagating the domain-wall, group 2 is in quiescent mode. After signal propagation is complete, group 2 devices go into signal propagation mode and stage 1 devices enter the quiescent mode to complete the entire signal propagation process. Using the above two methods, the interference between magnetic fringing fields are sufficiently suppressed for reliable domain-wall communications.

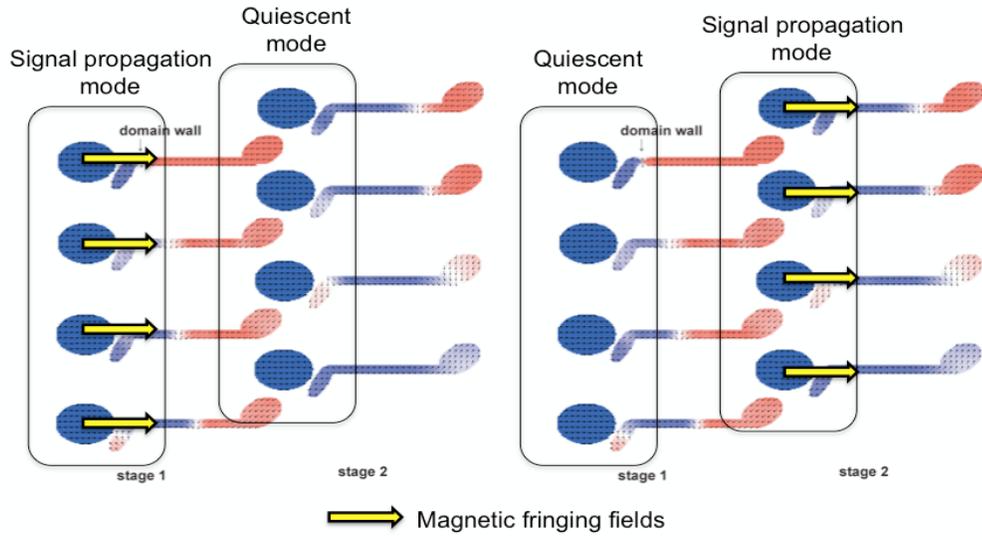


Figure 4.9. Two-stage domain-wall motion interconnect scheme

References

1. J. Inoue et al., “Nanomagnetism and Spintronics,” DOI: 10.1016/B978-0-444-53114-8.00002-9, Elsevier B.V. (2009)
2. V. Stojanovic, V. Oklobdzija, “Comparative Analysis of Master-Slave Latches and Flip-Flops for High-Performance and Low-Power Systems,” *IEEE J. Solid-State Circuits*, 35, No. 6 p876-884 (2000)
3. D. Atkinson *et al.*, “Magnetic domain-wall dynamics in a submicrometre ferromagnetic structure,” *Nature Mater.* **2**, 85–87 (2003)
4. S.S.P. Parkin, M. Hayashi, L. Thomas, “Magnetic Domain-Wall Racetrack Memory,” *Science* **320**, 190 (2008)
5. H. Meng, J.J. Wang, “A Composite free layer for high density magnetic random access memory with low writing field,” *IEEE International Magnetism Conference* (2006)

Chapter 5.

Applications of Magnetic Coupled Spin-Torque Devices

5.1 Introduction

Perhaps the greatest challenge facing any emerging logic device is the identification of specific applications where its strengths will be most keenly felt. There would be little advantage in constructing an entire microprocessor from magnetic logic elements; there may be great benefit to implementing a specific functional block within a hybrid system on a chip. Many people believe that the future of microelectronics lies in a diverse hybrid of technologies on a single platform, each doing what it does best [1]. In this chapter we will discuss three potential applications for MCSTD and investigate if MCSTD can augment CMOS to lead to a hybrid architecture that (hopefully) is faster, denser, lower power, and/or has increased functionality.

5.2 MCSTD based 2-bit full adder

5.2.1 Full adder design

As the first example application, we would like to examine the applicability of MCSTD gates to a general-purpose logic application: a full-adder. A Two-bit full-adder consists of combinational logic blocks that calculate the sum and carry-out. Sum and the carry-out can be expressed as follows:

- $\text{Sum} = A \oplus B \oplus C_{in}$
- $C_{out} = C_{in} \cdot (A+B) + A \cdot B = \text{majority}(A, B, C_{in})$

With MCSTD gates a two-bit full adder can be implemented as illustrated in Fig. 5.1. In overall, the compact and race condition-free MCSTD gate makes an adder circuit that is low power and highly reliable. Sum logic is cascaded MCSTD XOR gates. In MCSTD logic, XOR operation takes only one MCSTD gate to implement, which reduces the number of device count compared to any other logic family (Sec. 3.3.2). Two XOR gates are connected by a domain-wall motion interconnect (Sec. 4.4). Carry generation logic is expressed with a majority gate (see Fig. 5.1 (a)). A majority gate is available in MCSTD

similar to [2] because it can utilize anti-ferromagnetic coupling to manipulate its energy barrier. MCSTD majority gate does not have race condition errors, which is different from that of MQCA [2], All-Spin-Logic-Devices (ASLD) [3] and spin-wave device [19]. This is because MCSTD does not rely on placing its output device in a meta-stable state for logic operations, which all the aforementioned spin logic devices rely upon. At the same time, MCSTD gates use a synchronous circuit by combining clock signals with the switching current into the output device. In other words, MCSTD gates switch only when the switching current is applied to the output devices.

In this particular example, MCSTDs are assumed to be spin-valves. Since both spin-valve and MTJ emanate fringing magnetic fields, they can be used as a MCSTD gate. This is why we named our device magnetic coupled “Spin-torque device”. Spin-valve is considered here for low power consumption. Spin valves have an order of magnitude smaller absolute device resistance compared to MTJs, which can lead to lower power consumption. However, the switching current of the spin-valve is an order of magnitude higher than MTJ because of low spin-polarization. Having larger current results in greater power consumption although the resistance is low because the power consumption is quadratic with current. Our assumption is that current reduction in spin-valves is possible

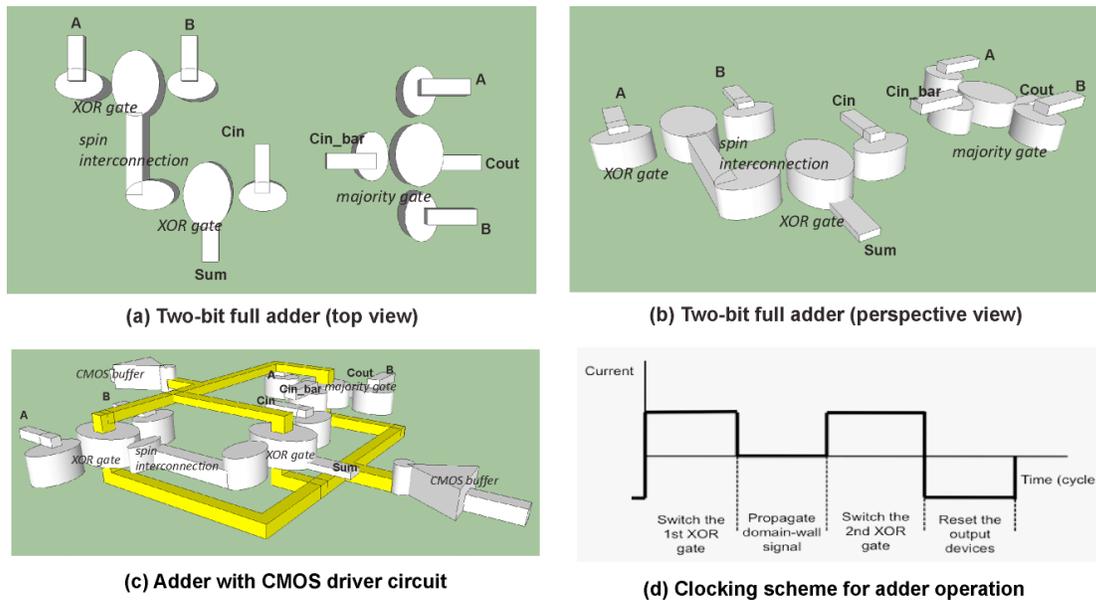


Figure 5.1. **Schematics of MCSTD 2-bit full-adder** (a) top view (b) three dimensional view (c) MCSTD full-adder with CMOS driver circuit and wires and (d) driver current waveform

through using materials, such as Heusler alloys, that intrinsically have near 100% spin-polarization [4]

Four clock cycles are needed for a MCSTD full-adder to perform an addition (Fig. 5.1 (d)). In the first cycle, switching current is applied to the first XOR gate and the switching is completed. Next, the XOR gate output signal propagates to the second XOR gate in the form of domain-wall motion. For reliable signal propagation, the output device in the second XOR enters the quiescent mode (see Sec. 4.4.2) to turn-off the fringing fields emanating from it. In the third cycle, switching current is applied to the second XOR gate and the majority gate for switching. After the computation result is collected at the output devices of the second XOR gates and the majority gate, all output devices are reset to ON state for the next operation. Note that the switching current polarity (or direction) has been changed for the reset (Fig. 5.1 (d)). MCSTD requires resetting the output devices because the logic operation in MCSTD depends on the initial state of the device, similar to CMOS dynamic circuits. The total count of spin-torque device switching operations is 6 because, there are three output MTJs and each switches twice including resets.

5.2.2 Energy estimate

As shown in Fig. 5.1, the two-bit full adder with MCSTD gates consists of two XOR gates and one majority gate. For each gate, there are two or three input devices and one output device. The input devices do not consume power because the signal propagation is

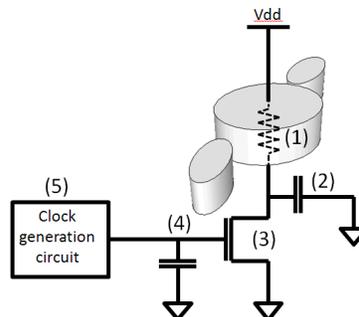


Figure 5.2. **Energy consumption components of MCSTD gates.** (1) switching energy of the output MTJ, (2) local interconnect capacitance (3) NMOS driver current source (4) global interconnect capacitance and (5) clock generation circuit

driven by the fringing fields of the output device. This makes the energy consumption of MCSTD gates equal to that of the output devices. The total energy consumption is divided into five components as denoted in Fig. 5.2 including that of the CMOS driver circuit.

The first energy component comes from the output spin-valve switching. We anticipate achieving $I_{\text{switching}}=30\mu\text{A}$, $R=10\Omega$ with a switching time of 1ns, which makes the energy consumption $I^2Rt=9\text{aJ}$ per switching. A recent work [5] demonstrated a spin-valve with $I=100\mu\text{A}$, $R=30\Omega$ with a switching time of 1ns. The output device reset doubles the energy consumption for the output device. With three output devices to switch, the total energy consumption in the MCSTD gates is

$$\begin{aligned} E_{MCSTD} &= (\text{No. of MCSTD gates}) \times (\text{switch \& reset}) \times (\text{energy per switching}) \\ &= 3 \times 2 \times 9 = 54 \text{ aJ/operation} \end{aligned}$$

The second component is the energy consumption in the wire that connects the output spin-valve to NMOS. This is greatly reduced by the low output voltage of the spin-valve, which is $V_{\text{out}}=I \cdot R=30\mu\text{A} \cdot 10\Omega=0.3\text{mV}$. Wire capacitance of $C=0.1 \text{ fF}/\mu\text{m}$ and wire length of $1\mu\text{m}$ are assumed.

$$E_{\text{local_wire}} = CV^2 = 0.1\text{fF} \times (0.3\text{mV})^2 = 9 \times 10^{-6} \text{ aJ/op}$$

Therefore, this second energy component is negligible.

The third component is the energy loss in the resistive MOS channel, $E_{RN\text{MOS}}$, equal to $I^2R_{\text{NMOS}}t$, where R_{NMOS} is the on-resistance of the MOSFET (in the triode region). Assuming $R_{\text{NMOS}} = R/3$, the energy becomes equal to 1/3 of the above-calculated value for the loss in the MCSTD device, i.e. 24.6aJ. Reducing R_{on} further would lower this energy component, but this comes at the expense of clocking energy, which is considered next.

The fourth component is the clocking energy consumption in the NMOS gate driver circuitry, which drives the NMOS gate and the wiring capacitance. This energy can be reduced substantially by employing a resonant clock generator scheme with charge recovery. If $R_{\text{NMOS}}=R/3$ and the charge recovery efficiency, $\eta=95\%$ are assumed,

$$\begin{aligned}
E_{NMOS} &= (1 - \eta) \times C_{gate} \times (V_{gs} - V_t)^2 = 0.05 \times 0.5 \times (V_{gs} - V_t) \times \frac{L^2}{(R_{NMOS} \times \mu_n)} \\
&= 0.05 \times 0.5 \times 0.5 \times \frac{(15nm)^2}{(10\Omega \times 400 cm^2/V \cdot sec)} = 6aJ/op
\end{aligned}$$

Energy consumption in the global wire from clock generator to NMOS gate can be amortized across the chip. In addition, the energy consumption is reduced to 5% due to the resonant charge recovery scheme described earlier. If we assume the per gate amortized global wire length to be 2 μ m,

$$E_{global_wire} = (1 - \eta) \times C_{wire} \times V^2 = 0.05 \times 0.2fF \times (0.5V)^2 = 2.5aJ$$

The fifth energy component, due to auxiliary circuitry in the clock generation can be amortized similarly, and will not affect the total energy (of a large VLSI chip) significantly. As a result, the total energy consumption in our MCSTD two-bit adder is approximately

$$\begin{aligned}
E_{total} &= E_{MCSTD} + E_{RN MOS} + E_{NMOS} + E_{global_wire} \\
&= 54 + 18 + 6 + 2.5 = 80.5 aJ/op
\end{aligned}$$

From the above calculation, two very important points need to be mentioned. First, we had to assume significant improvements over the current state-of-the-art spin-torque devices to achieve energy consumption comparable to that of CMOS. This result suggests further improvement is needed in spin-torque devices before they can be used as a logic device that can supplant CMOS.

Second, a considerable portion of energy is consumed in the CMOS circuit. A common problem of spintronic and non-volatile logic is that the switching of magnetic devices is “assisted” by a CMOS circuit. Other nanomagnet logic devices, such as MQCAs and ASLD, use initialization of magnetic devices into metastable states to make the switching easier. Our approach, MCSTD also requires reset of the output devices because the logic operation in MCSTD depends on the initial state of the device, similar to CMOS dynamic circuits. Placing magnetic devices in a particular magnetization or energy state is difficult to achieve with magnetic devices alone. ASLD and MCSTD use a “biasing” voltage. In order to minimize leakage power and to reverse the biasing

direction, the biasing voltage has to be driven by an active CMOS circuit. For spin-wave devices, external magnetic fields are necessary to create a resonant field. Most importantly, spin-wave devices have small output voltages and therefore require signal amplification. For logic operations, phase-matching and amplitude-equalization among spin-waves are needed and they all require CMOS circuits.

The consequence of having a CMOS circuit is that it limits the power consumption scaling. Let's assume 20% of my circuit needs to be CMOS gates. According to Amdahl's law in computer architecture, even if our new device consumes 1000 time less energy than CMOS gates, the maximum energy saving of the entire circuit is less than five.

$$Total\ energy\ saving = \frac{1}{0.2 + \frac{0.8}{1000}} \sim \frac{1}{0.2} = 5$$

The most desirable outcome for non-volatile logic is to come up with a logic device that can function without (energy inefficient) MOSFETs. If that is not possible, one alternative is to at least reduce the device count. Our approach, MCSTD is particularly efficient in reducing the device count. Having fewer devices helps to reduce 1) CMOS interconnection length and associated energy consumption 2) the number of NMOS drivers circuits for spin-torque devices and 3) local spin interconnection length. All spin devices use an interconnection mechanism that doesn't scale well with distance: spin-waves and spin-diffusion have limited traveling distance and longer MQCA chain of dots are more error-prone.

5.3 Logic embedded bio/image sensor

MCSTD holds great potential not only for generic logic, but also for logic-embedded sensors. This sensor can be in a far more efficient way than the current state-of-the-art in CMOS or by any other approach to non-volatile logic. From the multi-dimensional design space of MCSTD gates, new ways of building complex functions are possible, for example, XOR logic can be performed with a single MCSTD gate (see Chap. 3), which is

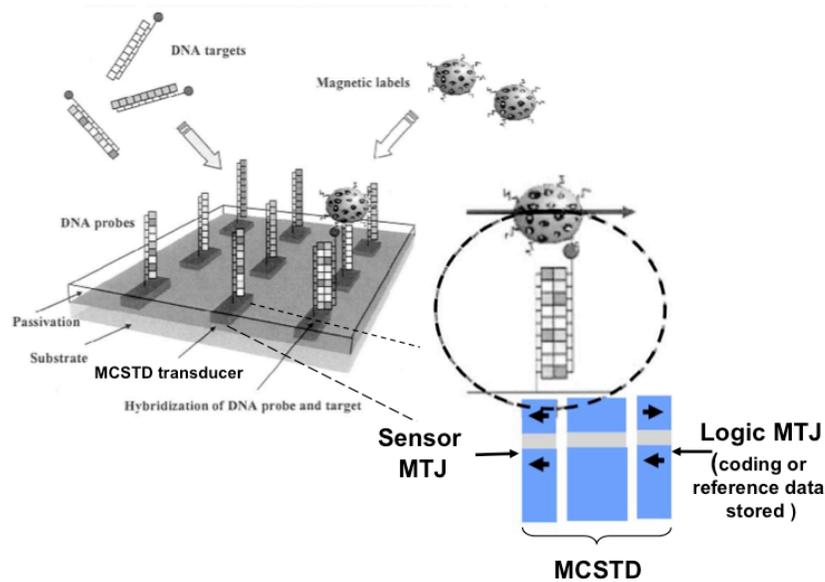


Figure 5.3. **MCSTD based logic embedded DNA microarray application**

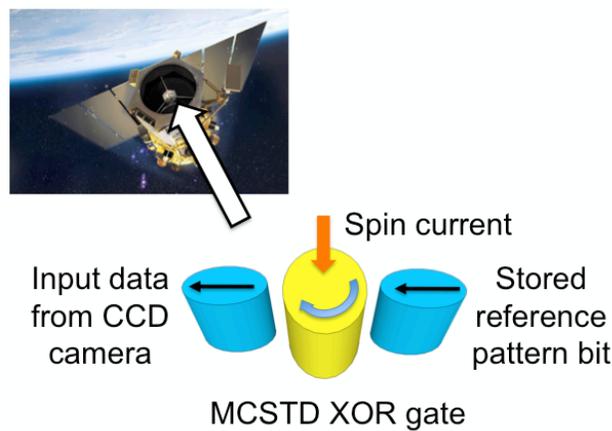


Figure 5.4. **MCSTD gates as a “smart” image sensor for surveillance applications**

more power and area efficient than in CMOS. Efficient XOR gate design together with its non-volatility enables MCSTD gates to address important sensor applications that require heavy signal processing and storage capacity. Such examples include DNA microarray (Fig. 5.3) and smart surveillance applications [6] as illustrated in Fig. 5.4. Because of the non-volatility of MCSTD gates, one can store an image or pattern and then instantly compare every incoming image with a simple array of MCSTD XOR gates without having to access memory as in conventional CMOS based schemes.

5.4 MCSTD for future reconfigurable logic

5.4.1 MCSTD in crossbar array architecture

In recent years, there has been a wide range of research efforts to find an optimum nano-system architecture for high-density future logic architecture. These architectures make use of nano-materials, such as Carbon NanoTube (CNT) [7], Rotaxane [8], nanowires [9,10,11] and NEMS switch [12]. General approach is to use nanowires or CNTs as interconnects to construct high-density crossbar array architecture (Fig. 5.5). Two-terminal programmable diode elements are placed at crosspoints to be used in logic operations. In addition, molecular-scale memory elements can be integrated at crosspoints to make the circuit non-volatile. High-density (due to crossbar array) and non-volatility makes this family of logic an excellent candidate for future reconfigurable logic.

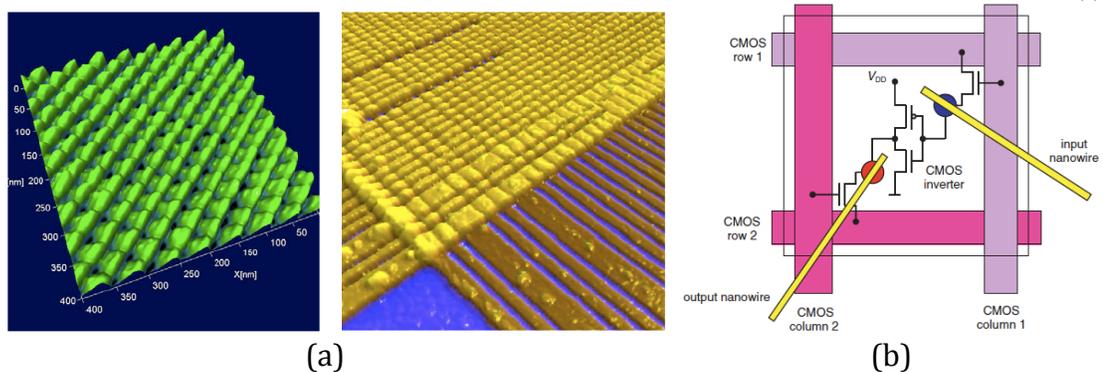


Figure 5.5. (a) Atomic force microscope topographs of a nearly defect-free region (left) and highly defective region (right) in a 34 nm pitch nanowire crossbar [11] (b) schematic of CMOL FPGA [10]

Although these previous works differ in terms of nanomaterials used, they all have the same issue of the limitations in diode-based wired logic: diode-like rectifying behavior in these two-terminal devices do not have a signal gain. Proposed solutions include 1) using nanowire FET type devices together with wired-diode based logic, 2) nano-architecture based on carbon nanotube (CNT) and CNT based MOSFETs and 3) using separate layers for different nano-materials and functionality: bottom layer is for Si CMOS and top layer is for nanowire-based memory/routing elements, post-silicon block

memories and CNT based vias that connects two layers [7]. Similarly, CMOL (CMOS/molecular hybrid) [10] proposes moving the inversion, gain part of logic into the CMOS layer, using the nanowires and crosspoint junctions only for wired-OR logic and signal routing. Another problem in a crossbar array with nano-devices is that additional CMOS gates are required for signal restoration. Two input nano-devices mentioned above typically have small R_{on}/R_{off} ratio and output voltage does not reach V_{dd} by device themselves.

As a result, challenges in crossbar array + nano-device architecture can be summarized as following: *high-density crossbar array architecture requires two-terminal logic devices but two-terminal devices (e.g., diodes) are not appropriate for logic device.* We believe this dilemma can be solved with the magnetic coupled spin-torque devices (MCSTDs) (Fig. 5.6, 5.7). MCSTDs consist of two terminal devices that fit very well in crossbar array architecture. On the other hand, MCSTD gates possess signal gain and other properties essential for logic applications (see Chap. 3). This is because, magnetic coupling between the input and output MTJs provides an additional third terminal that is needed to “gate” device operations similar to MOSFETs or BJTs. In other words, a MCSTD gate require only two terminals (top and bottom) for interconnection but, operates as three-terminal device. This advantage is very useful in building high-density

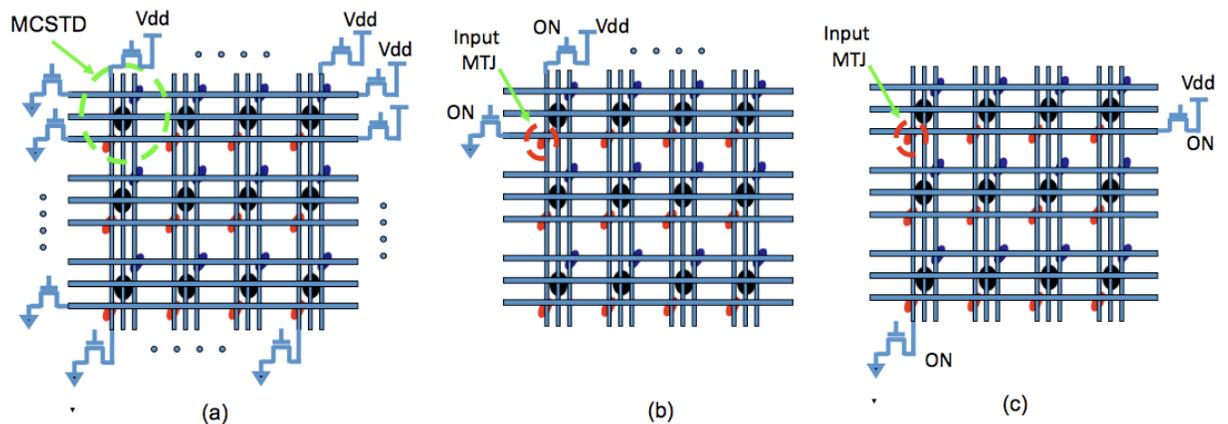


Figure 5.6. **MCSTD in crossbar array architecture.** (a) NMOS transistors (either connected to Vdd or GND) are wired to both ends of interconnects. Four transistors used to steer switching current direction into a single MTJ. For example, when two NMOS transistors are ON and the rest are OFF as shown in the figure, the current into the input MTJ (in red circle) flows (b) from bottom to top (c) from top to bottom contacts. NMOS transistors are shared by MCSTD gate in the same row or column.

logic device that can take advantage of crossbar array architecture. Finally, there is no need to restore output voltage to V_{dd} for interconnection because MCSTDs can communicate with the spin-interconnection or domain-wall motion interconnect scheme explained in Chapter 4. These interconnection schemes use current amplitude difference or magnetization as state variables inside the interconnection instead of fixed voltage levels of V_{dd} and GND. Otherwise, CMOS voltage amplifier are needed at the output of every spin devices, which is absolutely unacceptable overhead.

Operation of MCSTD gates in crossbar array architecture is explained in Fig. 5.7. The input MTJs are switched prior to the output MTJ. For switching, top and bottom wire that are connected to MTJs should be appropriately biased. As discussed in Chap. 2, the switching direction of spin torque device depends on the direction of current passing through the device. By appropriately biasing the top and bottom wires, the input MTJs get their magnetizations set. As shown in Fig. 5.6 (a), NMOS transistors (either connected to V_{dd} or GND) are connected to both ends of interconnects. Four transistors are used to steer switching current direction into a single MTJ. For example, when two NMOS transistors are ON and the rest are OFF (and left floating) as shown in Fig. 5.6 (b), the current into the input MTJ in red circle flows from bottom to top (Fig. 5.7 (c)). In contrast, if NMOS in Fig. 5.6 (c) and Fig. 5.7 (b) are turned ON and the rest are OFF, current flows from top to bottom contacts (Fig. 5.7 (b)). These NMOS transistors are shared by MCSTD gate in the same row or column. After the magnetizations of the input MTJs are set, current is sent through the output MTJ and it will switch or not switch depending on the biasing from the input MTJs.

5.4.2 Speedup of MCSTD over conventional MRAM technology

As seen in Sec.5.4.1, there are active research efforts to leverage the non-volatility of emerging future memories to build hybrid circuit with CMOS [13,14,15]. Among many emerging memory technologies, MRAM shows the fastest switching (write) speed and the longest endurance time [18]. Spin-torque Transfer MRAM, an MRAM technology that uses spin-torque transfer effect to rotate the magnetization instead of using a magnetic field, has a comparable speed to SRAMs (~2ns). Although MRAMs achieved tremendous speedup over conventional non-volatile memories, they are several orders of

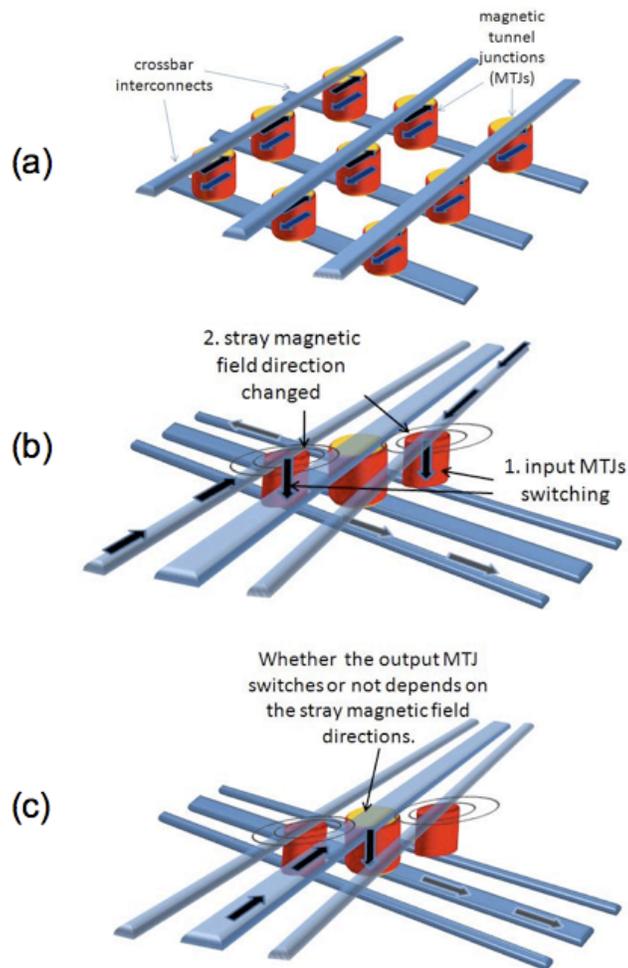


Figure 5.7. **MCSTD switching in crossbar circuit configuration.** (a) Typical crossbar circuit with two terminal devices at junctions. (b) MCSTD in crossbar generates a device gain (only one MCSTD is shown for brevity.) (c) output MTJ switching.

magnitude slower than CMOS logic devices. Without closing the device speed gap between MRAM technology and CMOS logic devices, putting them together to make hybrid circuits between MRAM and CMOS would not be feasible. We believe the MCSTD device has a better chance of closing the device speed gap with CMOS than any other emerging non-volatile device technology. This section discusses how MCSTD improves over MRAM switching speed and closes the speed gap with CMOS circuits.

a. Energy barrier lowering

Energy barrier modulation is the foundation of MCSTD logic and one can lower the energy barrier to reduce the write time of the non-volatile storage element. Switching

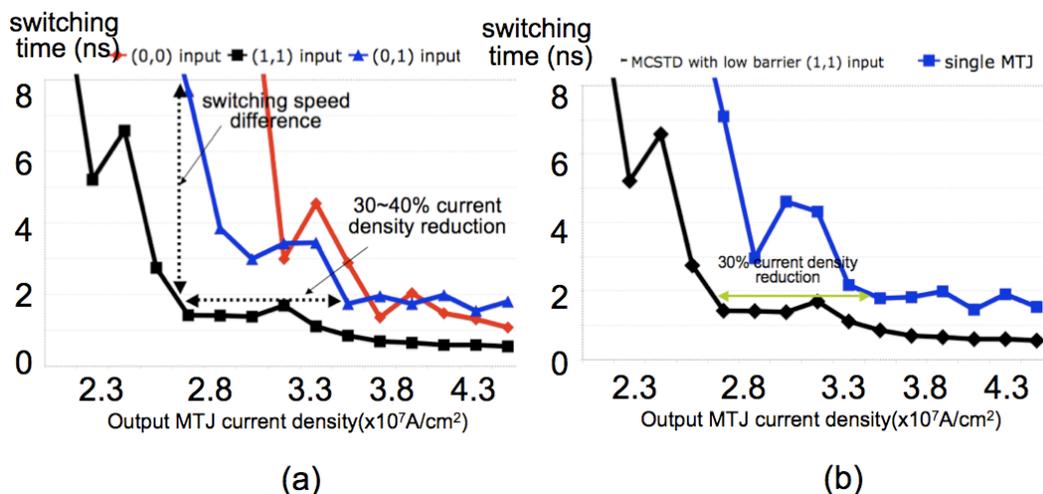


Figure 5.8. **Micromagnetic simulation of switching speed comparison** (a) MCSTD gate with different input signals (b) MCSTD versus single MTJ

speed of MCSTD depends on the magnetization of the input MTJs. Switching of MCSTD is fastest when its energy barrier is lowered by the input MTJs. Owing to the energy barrier lowering, the MCSTD switches two to four times faster than a single MTJ (Fig. 5.8).

b. Reduced device count

MCSTD devices can work as both logic and storage elements. If one builds a flip-flop storage with a NAND or NOR gate that precede or follow it (a logic embedded flip-flop), MCSTD circuit requires only one device. Compared to CMOS implementation, MCSTD circuit has reduced the number of logic stages and device count. Furthermore, by storing computational results in the device itself, communication delay and interconnection power consumption can be saved.

c. Heterogeneous material composition

A MCSTD consists of three parts: two input and output MTJs. Different materials can be used for each part to further improve the energy barrier modulation capability and reduce the write time (see Sec.3.3.5). If the output MTJ uses softer magnetic material than the input MTJ, the modulation in the energy barrier of the output MTJ becomes greater. To see the impact on the switching characteristics, the output MTJ free layer material is fixed as permalloy (NiFe) and the input MTJ free layer was varied with three

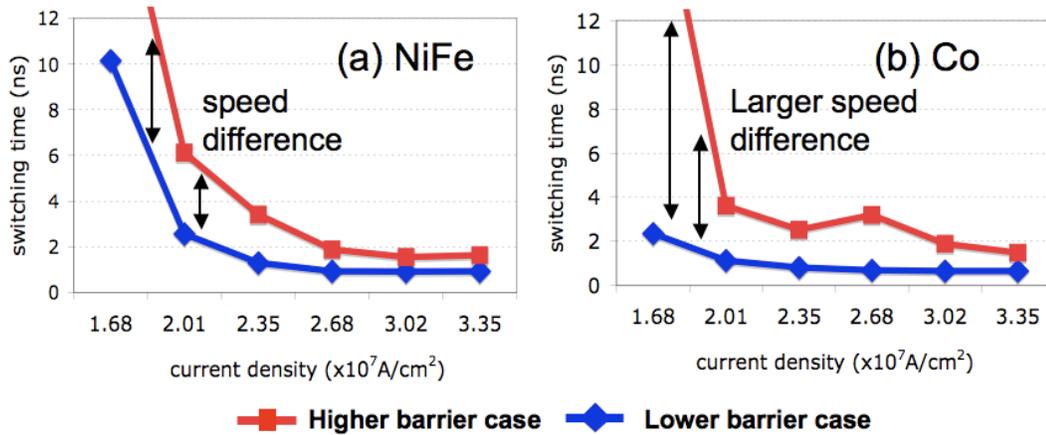


Figure 5.9. Switching speed comparison for different material used for the free layer of the input device. (a) NiFe has saturation magnetization of 800 G, (b) Co has 1200~1700 G. Co is harder material than NiFe.

different materials: NiFe, Co and CoFe. NiFe is the softest and CoFe is the hardest magnetic material. As the material for the input MTJ is changed to a harder material, the switching speed difference gets larger (Fig. 5.9). Using heterogeneous materials for the input and the output device is one effective way to increase the switching speed difference, which is favorable for the MCSTD logic device. One drawback is that the switching time of the input MTJ can increase.

d. Perpendicular Magnetic Anisotropy

Another improvement to reduce the MCSTD write time is to decrease the demagnetization field. The demagnetization field is a magnetic field inside of magnetic

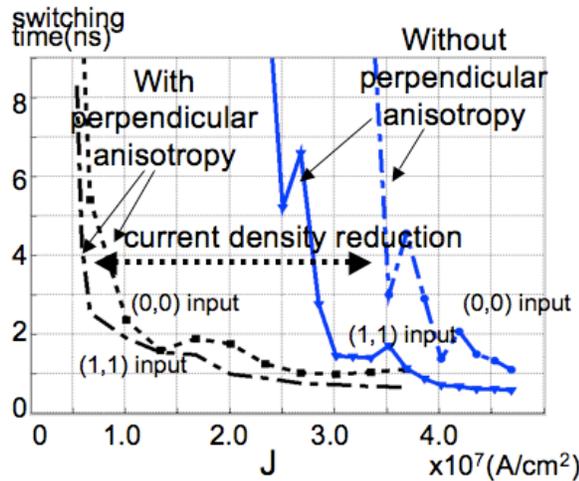


Figure. 5.10. Switching time improvement of MCSTD with perpendicular anisotropy

film that effectively suppresses the out-of-plane magnetization components. By reducing it, magnetizations in the MTJ free layer can easily be disturbed by a spin-torque current and switch to the other direction. By using material that has perpendicular magnetic anisotropy (PMA) of $K_{\perp}=1 \times 10^6 \text{ J/m}^3$, switching current density of a single MCSTD (sized $164 \times 264 \text{ nm}^2$) has been lowered by nearly an order of magnitude (Fig. 5.10). From the switching speed difference, we see that the energy barrier modulation of MCSTD is still effective although the difference is reduced. At higher PMA ($K_{\perp} > 2 \times 10^6 \text{ J/m}^3$), device shape can be circular instead of being elliptical. This is because all magnetizations are out-of-plane and there is no need for using the shape anisotropy of the ellipsoid to align magnetization. As a result, MCSTD devices size can be more compact and achieve high areal density.

5.4.3 Incorporation of MCSTD into CMOS

The ultimate goal of the CMOS/Nano-Magnetic hybrid circuit is to lower the two most significant barriers of modern computer systems: *communication barrier* and *power barrier*. First, we can greatly reduce the communication delay by merging logic and storage, which have been separated since the Von Neumann architecture [16]. Second, static power consumption can be nearly eliminated by making local storages non-volatile. As a first step to achieve this goal, we focus on Field Programmable Gate Arrays (FPGAs). FPGAs have SRAM based memory cells that need to be replaced by faster non-volatile memories. Also, most performance requirements for FPGAs used today are still low compared to ASICs or general purpose microprocessors, hence emerging devices slower than the state-of-the-art CMOS logic gates can be utilized.

For general purpose non-volatile logic, it is not necessary to make all logic gates non-volatile. If temporary storage elements in the datapath, such as flip-flops and latches are non-volatile, the wake-up time after power-gating will be no more than one clock cycle. This is because every pipeline stages have flip-flops and the non-volatile version of them can hold the data during power down period. This is a more power efficient way than to make entire datapath non-volatile. If all logic gates were non-volatile, they would maintain the states of the input or output ports which have to be cleared every clock cycle. This would result in additional overhead. Another reason for making flip-flops

non-volatile is the large percentage of clock time and power are already assigned to flip-flops in modern digital designs. Typically, flip-flops take up 20% of a clock cycle and a similar percentage of the power budget [17]. Emerging non-volatile devices are usually slower or more power hungry than CMOS gates. Thus, making only flip-flops non-volatile is a more practical approach than replacing all CMOS logic with non-volatile devices.

5.4.4 MCSTD based Look-up Table (LUT)

FPGAs are flexible system components because their logic functionality can be reconfigured. They “look up” the output for a logic computation in a look-up table (LUT) instead of having an assembly of logic gates for a computation. Logic functionality can be reconfigured simply by rewriting the contents of LUTs. LUTs are then connected together by an internal routing fabric in order to create larger computational blocks. LUT contents are stored in the memory cells. Usually, pass gate transistor based multiplexers

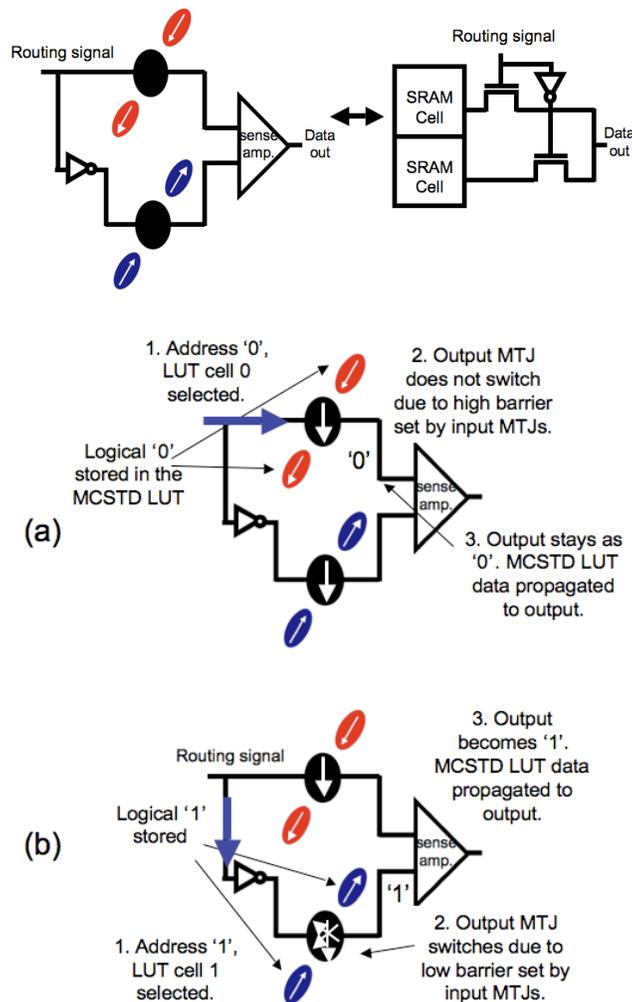


Fig. 5.11. Operation of MCSTD based FPGA Look-Up Table (LUT). One MCSTD replaces a SRAM cell and NMOS pass gate. (a) Retrieving data '0' (b) data '1' from MCSTD LUTs.

are used to retrieve look-up table values. The routing data that turn the pass gate transistors ON or OFF are stored in memory cells.

Currently, over 90% of the FPGAs use SRAM for LUT memory cell. Since SRAMs are volatile, the configuration state is lost whenever power is turned off. There has been effort to replace SRAMs with non-volatile memories such as Flash memory but has not been widely accepted in the industry. This is mainly because of the limited write counts of FLASH memory, for instance, up to 10^5 times (see Table 5.1). In addition, constant supply power is required for SRAM based FPGAs to maintain stored data in them, which becomes the source of static power consumption. Furthermore, NMOS pass gate in multiplexers have to be overpowered to transfer V_{dd} signal, which becomes another reason for power consumption. As a result, it is very desirable to have non-volatile LUT that has high endurance but does not need to be reprogrammed nor constantly connected to power supply.

MCSTD-based LUT achieves both non-volatility and device count reduction. First, spin-torque MRAM, the basic component of MCSTD, has higher endurance ($> 10^{15}$) than FLASH memories. Second, MCSTD reduces the device count by replacing CMOS multiplexor and SRAM cell. In detail, MCSTD can perform two-input logic operation such as NAND/NOR. Since multiplexers can be implemented with an AND gate,

	Field based MRAM (65nm)	Spin-torque MRAM (65nm)	FLASH (65nm)	DRAM (65nm)	SRAM (65nm)
Cell size (μm^2)	0.16	0.04	0.04	0.03	0.3
Read time (ns)	10	10	10~50	10	1
Write time (ns)	5	10	$10^5 \sim 10^8$	10	1
Program energy/bit (pJ)	100	1	10^4	5	5
Endurance	$>10^{15}$	$>10^{15}$	10^5 write	$>10^{15}$	$>10^{15}$
Non-volatility	Yes	Yes	Yes	No	No

Table 5.1. Comparison of advanced memories [18]

MCSTD can be used to build a 2:1 multiplexer. Two-input MTJs are tied together and store the data in look-up table. Address bit signal is connected to the output MTJ. When the address signal is '1', MCSTD is the same as the AND gate with one input pulled-up; the gate is transparent and the value in the input MTJs propagates to the output. If the address signal is '0' or V_{dd} , current flows into the output MTJ to drive it. If the signal is '0' or Gnd, no current will flow and the output MTJ will sit idle. As a result, one SRAM cell + NMOS pass gate can be replaced by a single MCSTD or two SRAM cells + one multiplexer can be replaced by a pair of MCSTDs and one sense amplifier can be shared by them as shown in Fig. 5.11. That is because, two neighboring memory cells do not get selected simultaneously and only one of them is read out. As a result, two neighboring MCSTDs are paired and the difference in their output signals is amplified for faster interfacing.

In summary, MCSTD based LUTs will eliminate tedious reprogramming and remove static power consumption due to its non-volatility.

5.4.5 MCSTD based Routing Fabric

Routing fabrics used in FPGAs need to have memory cells that maintain the routing information. In order to allow dynamic reconfiguration of routing fabrics, additional controllers or arbiters are required to process routing information. This is already a complex system because processing logic and storage elements should work together and stay updated. MCSTD gates can simplify much of this complexity (Fig. 5.12). NAND or NOR operations of MCSTDs can handle operations such as destination tag decoding of packet data. Non-volatility in MCSTD guarantees that the routing information is maintained even during a power glitch. Power can be aggressively turned off when there is no traffic in the routing, which will result in a huge saving in static power consumption.

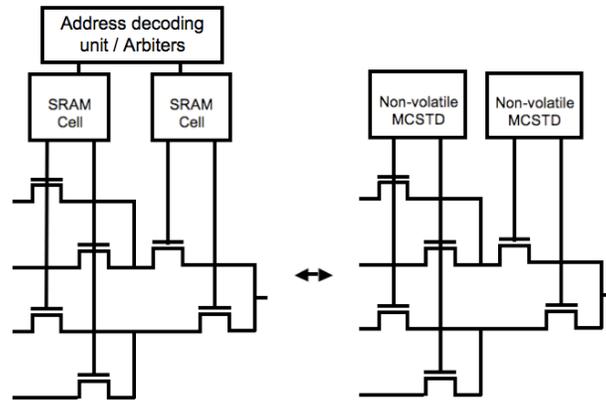


Fig. 5.12. MCSTD based reconfigurable routing fabric. SRAM cells and control units are replaced by MCSTD cells.

5.4.6 MCSTD based Flip-flop

MCSTDs make excellent non-volatile flip-flops as they work both as logic and storage elements. A clock signal is fed to the center output MTJ and the MCSTD will perform NAND/NOR/AND or OR logic function when the clock signal goes high. The output value is latched once the center MTJ has switched due to the non-volatility in ferromagnetic metal layers. Writing is much faster in MCSTD based flip-flop than other non-volatile flip-flops because, the energy barrier can be lowered significantly. During a read operation, a sense amplifier amplifies the differential TMR values between two complementary MCSTDs, which switch in the opposite direction at the same input. This eliminates the need for a reference signal, which is often used in other works [13,15]. MCSTD based non-volatile flip-flops will improve power/clock gating of CMOS logic to be much simpler and incur shorter wake-up times. Power gating transistors are not needed in MCSTDs because they don't require a connection to power supply other than the clock signal. Wake-up time is zero in MCSTDs as the data are stored in the device itself. In summary, MCSTD based flip-flop reduces the overhead of power-gating.

5. 5 Simulation result of MCSTD for future FPGA architecture

5.5.1 MCSTD vs. CMOS

To explore the benefits of magnetic coupled spin-torque devices (MCSTDs), first we compared their device area and power consumptions with CMOS gates. Please note that all results from simulations.

a. Area

Table 5.2 shows a comparison of MCSTD versus CMOS device counts for different logic functions. One of the unique advantages of MCSTD is that a single device can perform one logic operation: it takes only one MCSTD gate to do NAND, NOR, XOR and XNOR operations respectively (see chap. 3). For the spin interconnection scheme (Sec. 4.1), complementary gates are needed and that makes the required device count to double in Table 5.2. Note that no device is needed for an inverter in MCSTD because, all that is needed is electrical vias that can be placed in the empty spaces. This is still a far lower device count than that of CMOS. Next, Table 5.3 illustrates the actual layout size comparison with standard cell CMOS gate libraries. MCSTD shows three order of magnitude small gate sizes. In summary, the lower device count and smaller device area of MCSTD will allow it to be lower power than CMOS.

	NAND	NOR	XOR	XNOR	inverter
CMOS	4	4	16	16	2
MCSTD	2	2	2	2	0

Table 5.2. Device count comparison between CMOS and MCSTD. (CMOS: number of transistors, MCSTD: MCSTD gate count)

b. Power

Figure 5.13 (a) shows the power-delay product of MCSTD and CMOS two-input NOR gates in different technology nodes. For MCSTDs, switching current density becomes the key metric that determines the switching delay and power consumption of the device. For micro-magnetic simulation, current density of $J \sim 5 \times 10^5$ A/cm² is used. This is roughly half of the best switching current density reported for spin-torque transfer RAM (STT-RAM). Switching current density for MCSTD is lower because, 1) there is

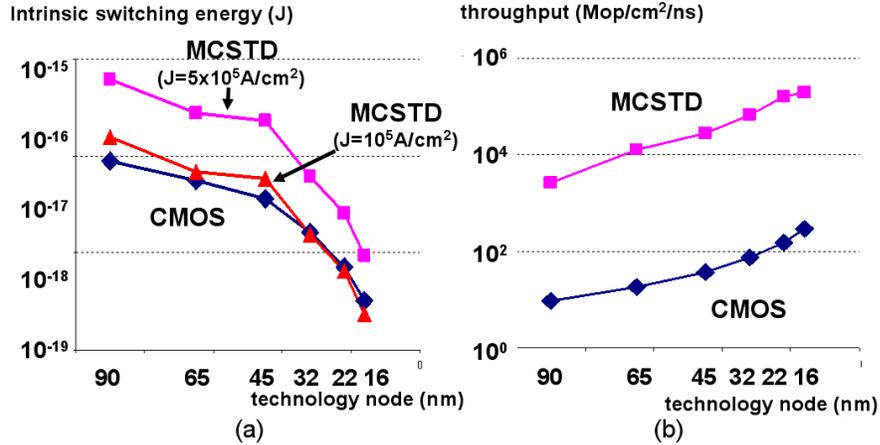


Figure 5.13. MCSTD intrinsic switching energy (a) and throughput (b) in comparison with CMOS.

energy barrier reduction due to the biasing fringing fields from the input MTJs and 2) data retention time requirement for magnetic logic device can be much smaller than memory device, which can be reduced from 10 years to 10 days roughly. Thus, thermal stability can be reduced to lower the switching current density in MCSTD gates. However, MCSTD consumes an order of magnitude larger energy than CMOS. The differences between CMOS and MCSTD are getting smaller as device scaling continues. If switching current density can be reduced to $J \sim 1 \times 10^5 \text{ A/cm}^2$, MCSTD gate consumes less energy than CMOS gates. This trend will become more obvious when we take the fact that MCSTD XOR/XNOR gates take fewer devices than CMOS and no device except electrical vias are needed for inverters. Figure 5.13 (b) compares the throughput of CMOS and MCSTD gates. MCSTD achieves several orders of magnitude larger throughput than CMOS because, more MCSTD devices can be packed into the same area despite the longer switching delay.

5.5.2 MCSTD based FPGA Look-up Table (LUT)

a. Area and transistor count

To implement a two-input LUT, nano-magnet spin-torque device based circuit requires four MCSTDs and one 2:1 address decoder. Each MCSTD stores configuration data in its input MTJs, which are non-volatile storage. On the other hand, CMOS FPGA LUT takes four SRAM cells and a 2:1 multiplexor. LUT does not need a sense amp. Also, depending

hSpice Model	V _{dd}	Delay	SRAM Power	SRAM leakage Power	2:1 mux Power	2:1 mux leakage Power	Total Energy Consumption
PTM 32nm	0.9V	196ps	14.7uW	0.02uW	2.8uW	110nW	3.45fJ
PTM 45nm	1.0V	200ps	22.6uW	5.54nW	5.4uW	45.5nW	16.3fJ
TSMC 0.18um	1.8V	334ps	126uW	49pW	66.8uW	0.19nW	64.4pJ
MCSTD (32nm width)	Supply current	MCSTD LUT delay	MCSTD address decoder delay				8.9fJ
	40uA	200ps	409ps				

Table 5.4. Comparison of CMOS LUT and nano-magnet/CMOS hybrid circuit LUT energy consumption

on the operating frequency of the LUT, a sense amplifier is reasonable so as not to pay a large delay penalty waiting for the bit lines to swing full rail. Thus, in future FPGA systems, sense amplifiers may be needed.

b. Energy consumption

Regarding energy consumption, a CMOS LUT in 0.18um, 45nm and 32nm technology nodes are compared with MCSTDs of 200nm, 100nm, 50nm and 32nm node (for example, 32nm node of MCSTD is defined as the average of output MTJ width and length is 32nm) (Table 5.4). Similar to CMOS device scaling, power consumption of spin-torque devices scales down with physical device scaling. This is because it takes less current to switch the magnetic devices. On the other hand, the resistance of the magnetic devices goes up when nano-magnetic devices are scaled down. Resistance increase of the MTJ is a two-edged sword: large resistance increases the power consumption but it makes the output signal large and the sense amplifier can catch the signal more quickly.

For the simulation, the switching of the MCSTD gates are simulated with micromagnetic simulator, OOMMF, and its results are fed into HSpice, where the MCSTDs are modeled as variable resistors. MTJs are assumed to have $1.0\Omega\cdot\mu\text{m}^2$ of resistance-area product and tunneling magneto-resistance (TMR) of 400%.

c. MCSTD based FPGA Look-up Table (LUT)

As described in Chapter 3, a MCSTD gate not only replaces the SRAM cell, it also provides a 2:1 multiplexer functionality by performing AND operation. Fig. 12 shows the waveforms of the signals when LUT values are accessed. For the simulation, the TMR change of the MCSTD gates are simulated with OOMMF and fed into hSpice, where the MCSTDs are modeled as variable resistors. MTJs are assumed to have $1.5\Omega\cdot\mu\text{m}^2$ of resistance-area product and spin current density used is $1.39 \times 10^6 \text{A}/\text{cm}^2$. As address bit 0 goes high, MCSTD-based LUT cell_1 gets selected and the output MTJ of it attempts switching. AND gate with one input 'high' is transparent to the other input, which is the data stored in the input MTJs. For LUT operations, two input MTJs are tied together and they stored the LUT cell value. In our simulation, the input MTJs are holding data '1' and it showed up at the output of sense amplifier after 0.65ns of delay (see Fig. 5.14). Notice

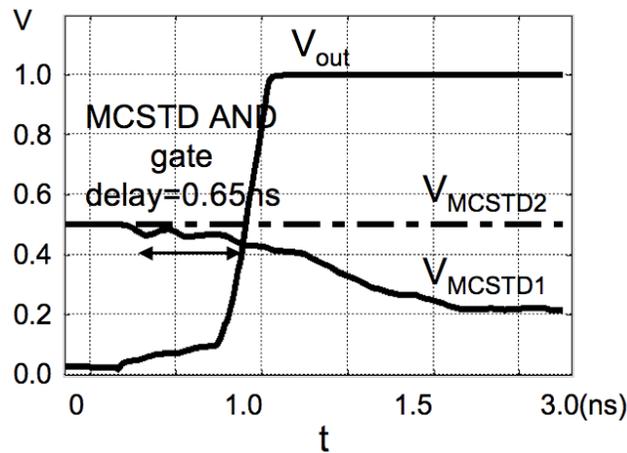


Fig. 5.14. MCSTD based FPGA LUT read operation. As MCSTD1 (LUT cell 0) switches to ON stage, V_{out} of the sense amplifier rises to V_{dd} . Gate delay is shorter than average SRAM read time.

that data '1' in the input MTJs lower the energy barrier height of the output MTJ and facilitates the switching. If a single MTJ or MRAM cell was used for this purpose, the switching time can be more than three times slower as there is no energy barrier lowering effect. MCSTD-based LUT access time is much faster than normal SRAM cell read time of 2 ns. By using a MCSTD gate, one sense amplifier and two pass gates for multiplexer could be saved, which is 11 transistors. More importantly, LUT cells don't need to be

constantly powered and overpowering NMOS pass gates with voltage higher than V_{dd} to compensate for poor transferring capability of passing logic '1'.

5.5.3 MCSTD based Routing Fabric

Fig. 5.15 shows the signal waveforms during dynamic reconfiguration of FPGA routing fabric. As an example, the MCSTD performs a NOR operation for address processing operations and then stores the routing information in the device itself. Total reconfiguration time is only 1.6ns. This delay is similar to an L2 cache miss in the state-of-the-art general purpose processors, which can bring “true” reconfigurable processors much closer.

MCSTD gates simplify FPGA routing fabric design by replacing a logic gate and one SRAM. Switching current of MTJs linearly decreases as the device dimension scales.

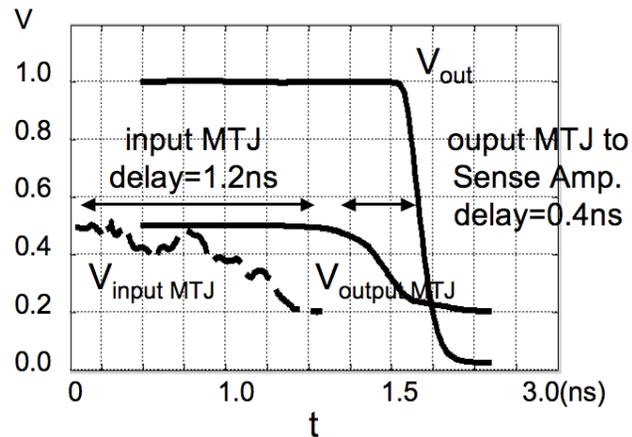


Fig. 5.15. MCSTD based routing fabric reconfiguration process. Input MTJs receive the input values first and the output MTJ switches to perform address processing related simple logic operations. Results show up at the sense amplifier. Delays for each stage are shown.

Due to high aspect ratio and the lack of energy barrier lowering makes it more difficult to switch the input MTJs than the output MTJs. For the simulation result shown in Fig. 5.15, current density of $J=9 \times 10^7 \text{A/cm}^2$ was used, which is much higher than that of the output MTJ ($1.39 \times 10^6 \text{A/cm}^2$).

5.5.4 MCSTD based Flip-flop

The MCSTD-based flip-flop was simulated for general purpose non-volatile logic. It receives two inputs and performs logic operations such as NAND/NOR/NOT. Fig. 5.16 shows MCSTD flip-flop waveforms performing NAND operation. At (0,0) input, the flip-flop delay is 0.35ns (Fig. 5.16 (c)) and the resulting power-delay product (PDP) is 15.5fJ for a single output MTJ. The PDP of the input MTJ is 17.4fJ. As MTJs scale down, magneto-resistance increases but the switching current decreases faster, which scales down the power consumption linearly. The PDP of the MCSTD based flip-flop under simulation is 72.3fJ, which is 0.7~3.2 times larger than 0.20um node technology CMOS flip-flops [16]. When compared with the 32nm node technology CMOS flip-flops, the PDP of MCSTD-based flip-flop is roughly ten times larger. This is mainly because of the high switching current density and high TMR of MTJ device. When the static power consumption is included in the comparison, the total PDP of CMOS flip-flop breaks even at 6.8ns becomes larger onwards. Usual stand-by time of embedded systems is much longer than this breakeven point and MCSTD-based flip-flop can be considered as more energy-efficient for embedded devices with long wait time.

MCSTD-based flip-flop gate delays are different for different input values. This is

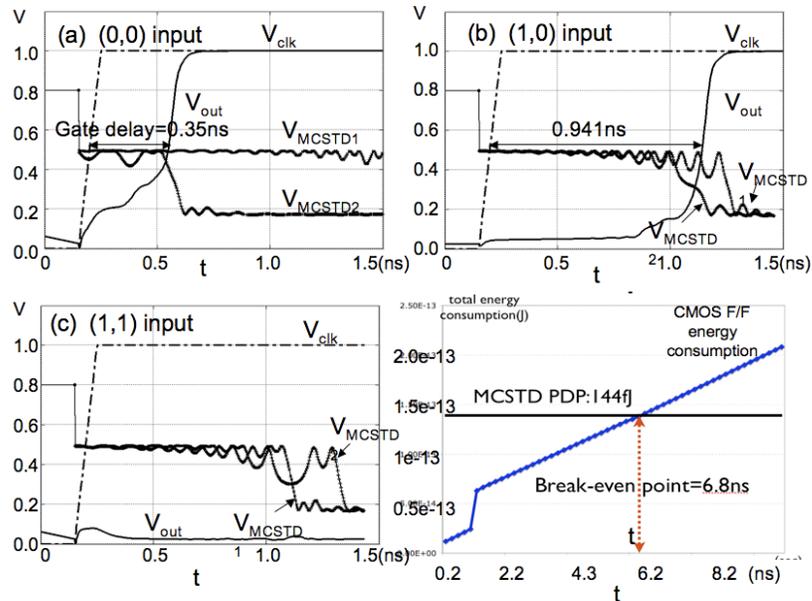


Fig. 5.16. Voltage waveforms of MCSTD based Flip-flop. (a) (0,0) (b) (1,0) (c) (1,1) input (d) Switching time improvement of MCSTD with perpendicular anisotropy.

because the logic operation embedded in the MCSTD-based flip-flop incurs different delay. The delay is the largest when the input is (0,1) or (1,0) (Fig. 5.16(b)). For these inputs, stray fields from the input MTJs are opposing each other and cancel. Thus, there is no energy barrier lowering to help the output MTJs to switch faster.

References

1. R. P. Cowburn, M. E. Welland, "Room temperature magnetic quantum cellular automata," *Science* 287, pp.1466–1468 (2000)
2. A. Imre et al., "Majority logic gate for magnetic quantum-dot cellular automata," *Science* 311, pp.205–208 (2006)
3. B. Behin-Aein, D. Datta, S. Salahuddin, S. Datta, "Proposal for an all-spin logic device with built-in memory," *Nature Nanotechnology*, DOI: 10.1038 (2010)
4. T. Kubota *et al.*, "Half-metallicity and Gilbert damping constant in $\text{Co}_2\text{Fe}_x\text{Mn}_{1-x}\text{Si}$ Heusler alloys depending on the film composition," *Appl. Phys. Lett.* 94, 122504-122503 (2009)
5. S. Mangin, Y et al., "Reducing the critical current for spin-transfer switching of perpendicularly magnetized nanomagnets," *Appl. Phys. Lett.* vol. 94, 012502-012503 (2009)
6. L. Liu, S. Kesavarapu, J. Connell, A. Jagmohan, L. Leem, B. Paulovicks, V. Sheinin, L. Tang, H. Yeo, "Video Analysis and Compression on the STI Cell Broadband Engine Processor", *Proc. International Conference on Multimedia and Expo* (2006)
7. C. Dong, S. Chilstedt, D. Chen, "Reconfigurable circuit design with nanomaterials," *Proc. Design Automation, Test in Europe* (2009)
8. Y. Chen et al, "Nanoscale molecular-switch crossbar circuits," *Nanotechnology* 14 pp.462~468 (2003)
9. A. DeHon, M. J. Wilson, "Nanowirebased sublithographic programmable logic arrays," *Proc. International Symposium on FPGAs*, pp.123-132 (2004)
10. D. Strukov, K. Likharev, "CMOL FPGA: a reconfigurable architecture for hybrid digital circuits with two-terminal nanodevices", *Nanotechnology* 16 (2005)
11. G. S. Snider, R. S. Williams, "Nano/CMOS architectures using a field-programmable nanowire interconnect", *Nanotechnology* 18 (2007)
12. Y. Zhou et al., "Low Power FPGA Design Using Hybrid CMOS-NEMS Approach", *Proc. International Symposium on Low Power Electronics and Design*, 2007

13. A. Mochizuki, H. Kimura, M. Ibuki and T. Hanyu, "TMR-Based Logic-in-Memory Circuit for Low-Power VLSI," *IEICE Trans. Fundamentals*, vol. E88-A, No.6, pp.1408-1415 (2005)
14. W. Wang, A. Gibby, Z. Wang, T. Chen, S. Fujita, P. Griffin, Y. Nishi, S. Wong, "Nonvolatile SRAM Cell," *Proc. International Electron Device Meeting*, DOI: 10.1109/IEDM.2006.346730 (2006)
15. W. Zhao, E. Belhaire, B. Dieny, G. Prenat, C. Chappert, "TAS-MRAM based Non-volatile FPGA logic circuit," *Proc. International Conference on Field-Programmable Technology* (2007)
16. V. Stojanovic, V. Oklobdzija, "Comparative Analysis of Master-Slave Latches and Flip-Flops for High-Performance and Low-Power Systems", *IEEE J. Solid-State Circuits*, 35, No. 6 pp. 876-884 (2000)
17. http://en.wikipedia.org/wiki/Von_Neumann_architecture
18. D. Bondurant, B. Engel, J. Slaughter, "MRAM- The future of non-volatile memory?" *Portable Design*, July (2008)
19. A. Khitun, M. Bao, J.Y. Kim, A. Hong, A.P. Jacob, K.L. Wang, "Logic devices with spin wave buss: potential advantages and shortcoming," *Proc. Device Research Conf.*, 2008

Chapter 6.

Fabrication process of Magnetic Coupled Spin-Torque Devices

6.1 Fabrication challenges

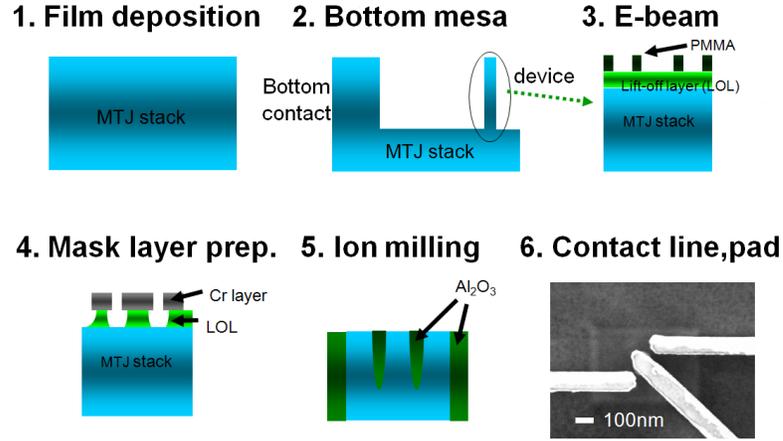


Figure 6.1. Fabrication process flow of MCSTD gates

There are several challenges to the fabrication of MCSTD gates. These include: (1) patterning small area MTJs which are sufficiently close together (<30 nm) for strong magnetic coupling, but remain electrically isolated, (2) ion-milling the closely spaced MTJs whilst maintaining sidewalls of high integrity, (3) making electrical contacts to each of the closely spaced MTJs. In this chapter, we describe the fabrication process of MCSTD gates (Fig. 6.1) and our efforts to overcome these fabrication challenges.

The fabrication process of MCSTD gate discussed in this chapter was developed and performed by Stuart S. P. Parkin group at IBM Almaden Research Center. I participated as a project member of this fabrication process and provided device dimensions and feedbacks from the designer's point of view. I gratefully acknowledge Parkin group for providing the state-of-the-art facilities and long years of expertise in MTJ fabrication, which was indispensable element in the successful fabrication of MCSTD gates. MTJ film stack was prepared by Dr. See-hun Yang, e-beam lithography and Cr mask preparation by Dr. Charles Rettner and ion-milling by Mr. Brian Hughes and probe-station prepared by Dr. Xin Jiang.

6.2 MTJ film stack

First, the MTJ film stack is deposited by magnetron and ion beam sputtering. All the MTJ films are deposited on thermally oxidized 1-inch Si wafers. A small in-plane magnetic field was applied during deposition to define the easy axis of the magnetic films. The MTJ film structures used are (from bottom to top),

100Å Ta | 300Å Ir₂₂Mn₇₈ | 6Å Co₄₀Fe₄₀B₂₀ | 30Å Co₇₀Fe₃₀ | 8Å Ru | 27Å Co₇₀Fe₃₀ | 8Å Mg | 4Å Mg in (95 Ar/5 O₂) | 20Å Co₄₀Fe₄₀B₂₀ | 50Å Ta | 50Å Ru

100Å Ta | 300Å Ir₂₂Mn₇₈ | 6Å Co₅₆Fe₂₄B₂₀ | 24Å Co₇₀Fe₃₀ | 5Å Ru | 27Å Co₇₀Fe₃₀ | 8Å Mg | 3Å Mg in (95 Ar/5 O₂) | 20Å Co₅₆Fe₂₄B₂₀ | 50Å Ta | 50Å Ru

MgO thickness and CoFeB (free layer) thickness and composition were varied to adjust the required switching voltage and the coercivity of the output MTJ.

The fixed layer of the MTJ consists of an IrMn exchange bias layer and a CoFeB|CoFe|Ru|CoFe Synthetic Anti-ferromagnetic layer (SAF) [1]. Exchange bias from underlayer (IrMn in our case) "pins" the fixed layer. After MTJ film stack deposition, the devices were annealed at 260~300°C for 30 minutes in a 1T (=10⁴ Oe) field. This is to increase the exchange bias from IrMn layer.

Synthetic Anti-Ferromagnetic layer forms a closed-loop between the layer above and below Ru. The purpose of the closed loop is to have

- 1) No interlayer coupling with free layer. If there is no magnetic flux going outside the SAF layer, there will be no effect on the free layer.
- 2) SAF layer is seen non-magnetic from outside. When magnetic field is applied from outside it doesn't feel any fields and have no rotation. This make the SAF layer perfectly pinned.

6.3 Bottom contact

As the first step of nanofabrication, a mesa structure is patterned. It is simply a rectangular area of MTJ stack unetched. The rest of region is etched down to Si. A

negative e-beam resist Hydrogen Silsesquioxane (HSQ) and Durimide (lift-off-layer, LOL, 200nm thickness) double layer to make mesa patterns.

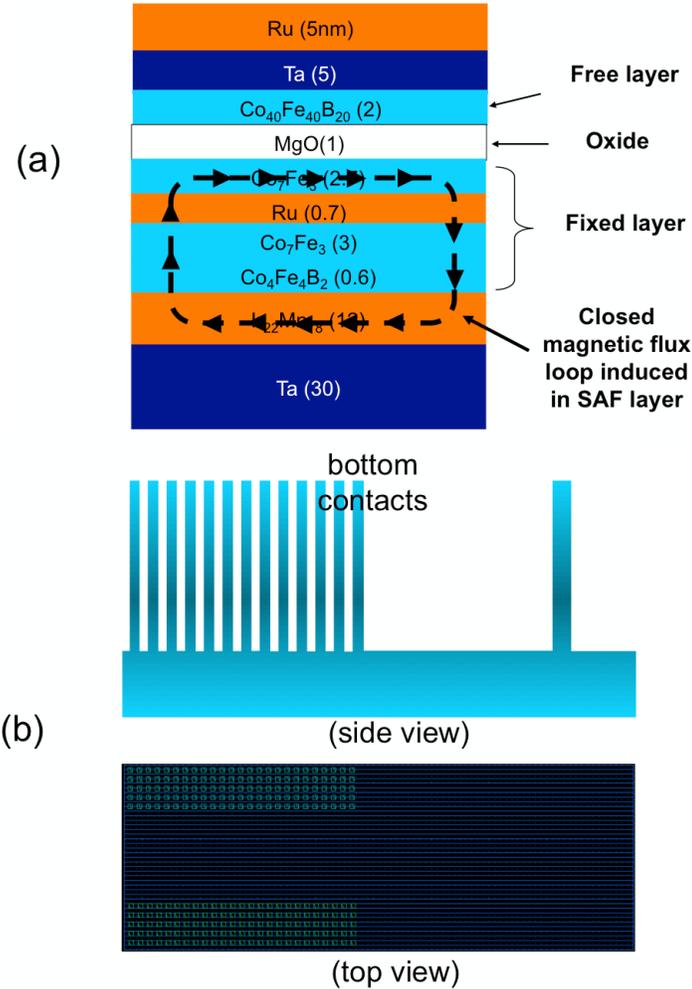


Figure 6.2. (a) MTJ film stack used for MTJ fabrication and (b) bottom contact structure (side and top view)

For bottom contact, a group of nano-pillars (125 square pillars each sized $1\mu\text{m}^2$) are patterned and etched with ion-milling technique. Since the voltage drop across this group of nano-pillars is negligible, the voltage level at the bottom contact can be read at the top of the nano-pillars. Resistance of bottom and top electrical contact should be carefully estimated and minimized. This is because,

$$R_{Total} = R_{Bottom_contact} + R_{MTJ} + R_{Top_contact}$$

Only MTJ resistance is magneto-resistance (MR) and changes. If top and bottom contact

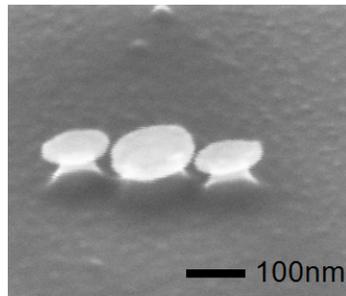
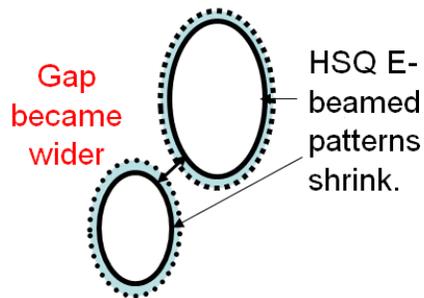
resistance are overwhelming, MR change will not show up in the total resistance measurement.

When removing HSQ resist layer, it is important to have alumina thin (alumina deposited after ion-milling to protect the MTJ film stack sidewalls). Otherwise, resist strippers such as N-methyl 2 pyrrolidone (NMP) cannot attack and dissolve the under-layer resist. NMP never dissolves alumina. If there's alumina blocking NMP meeting the under-layer resist, liftoff will never work. To break away alumina sidewalls, process called “Snow-clean” is used. It is to apply strong flow CO₂ gas to the surface of the wafer followed by heating the wafer at 120°F. It makes crack on Al₂O₃. Dip the wafer in hot NMP (or acetone) with ultrasonic.

6.4 E-beam lithography

Generally, negative e-beam resists, such as Hydrogen Silsesquioxane (HSQ) are used for nano-pillar e-beam lithography. However, HSQ was avoided in this work because of potential undesirable shrinkage, which will increase the critical element spacing. All of the following conditions can contribute to HSQ resist shrinkage: (1) time delay and ambient conditions between the e-beam exposure and subsequent development, (2) chemical cross-linking between HSQ and the underlying resist layer, (3) development time if the insoluble layer is removed by hydrofluoric (HF) acid, and, (4) long exposure time [2, 3]. Shrinkage of the HSQ resist is often favored to produce smaller dimension than what can be defined with e-beam lithography. However, HSQ shrinkage in this case could lead to loss of control on the spacing between the MTJs (Fig. 6.3). The spacing between the MTJ devices is a critical dimension of the MCSTD, which can affect the dipole coupling strength between these MTJs. In addition, fine details of the device patterns can be lost due to HSQ shrinkage. These effects will likely be more important when the MCSTD devices are scaled to smaller dimensions.

To produce the effect of negative resist without using HSQ, a positive resist, such as PMMA, was used as the e-beam resist and a Cr hard mask layer was prepared on top of each MTJs (Fig. 6.3). First, spin-coat PMMA (resist) and Durimide (LOL) double layer. Then do e-beam lithography. After exposure, PMMA was post-baked for 1 min at 120°C. Otherwise resist boundary will be rough. Post-baking hardens insoluble region less soluble and makes soluble region more soluble. Next, exposed pattern was developed with Methyl-Iso-Butyl-Ketone (MIBK). MIBK is used instead of anisole as a base for PMMA. This is because Durimide dissolves in anisole. For high resolution PMMA development, low temperature (near 0°C, 0.3°C) developer + ultrasonic were used.



Chrome + Lift-off layers

Figure 6.3. Fabrication issues of using negative e-beam resist for MCSTD fabrication (top) and Chrome hard mask layer as a solution (bottom)

6.5 Chrome hard mask layer

Do plasma etch to get PMMA patterns transferred to Durimide layer (Fig. 6.4). We used O₂ plasma etching but it has some undercut. Undercuts are good for lift-off but are not good for transferring patterns to bottom layers. CO₂ etching, which has less undercut can be an alternative. Deposit chrome (through e-beam evaporation) where Durimide is

exposed. Use acetone to etch remaining PMMAs. PMMAs are soluble to acetone but Durimide is not. N-methyl 2 pyrrolidone (MNP) can be used in place of acetone. Cr adheres well to Durimide. Adhesion was so strong that it stayed on the Durimide when ultrasonic cleaning was strong enough to destroy the Si wafer.

Appropriate thicknesses for the Cr mask layer and lift-off layer were found by trial-and-error: if the thickness of Cr mask layer is too thin (e.g., <7nm), alumina sidewall will become too thick (similar to HSQ removal in Sec. 6.3) and not allow the etchant to reach the lift-off layer. Lift-off resist layer thickness was limited not to affect the e-beam lithography resolution.

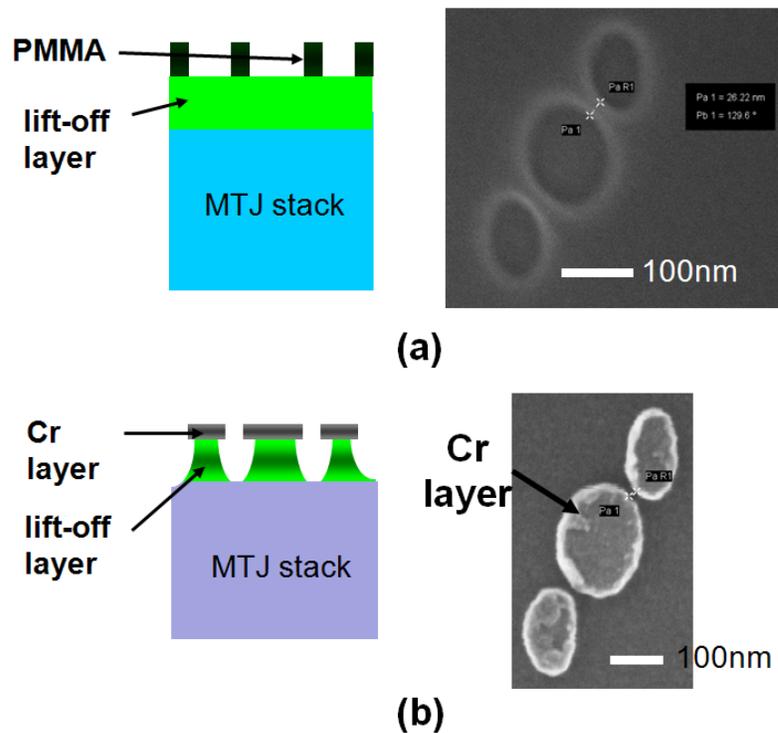


Figure 6.4. E-beam lithography and Cr hard mask layer preparation for ion-milling

6.6 Ion-milling the gaps between MTJs

A two-step ion-milling process was utilized to etch away the narrow gap between the input and output MTJs. First, do etching with ion-beam at 5° angle normal to the surface

until the MgO barrier is etched through. Then, two minute 80° angle (normal to the surface) ion-beam for cleaning the re-deposited material (Fig. 6.5). Any conducting path across the tunneling barrier formed from redeposition can severely degrade TMR of MTJs. Second, do etching with ion-beam at 5° until IrMn layer is reached (i.e., IrMn detected by in-situ Secondary Ion Mass Spectrometry (SIMS) analysis). This is followed by 8 minute 80° angle (normal to the surface) ion-beam for cleaning the re-deposition. Milling time and beam angles should be adjusted for better results.

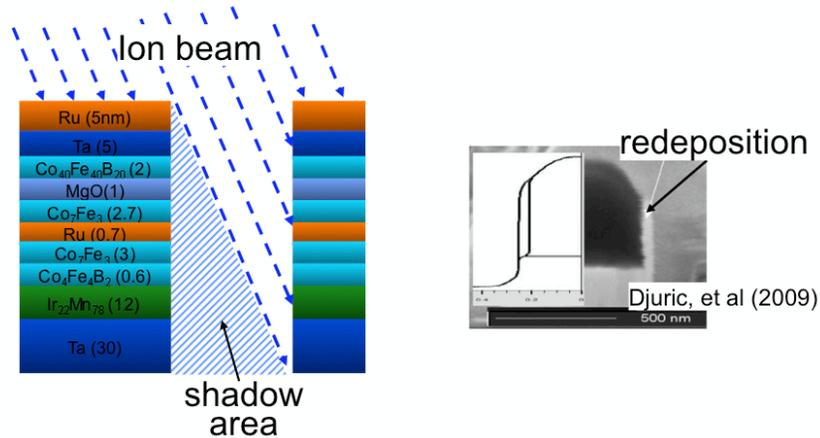


Figure 6.5. **Ion-milling process and shadow area.** Reducing the redeposition on the side walls of nanopillars is one of the key fabrication challenges

Deep and shallow ion-milling depths (of MTJs) have following trade-offs. If the etched depth is deep, lift-off process becomes difficult. This is because more “filler” material such as alumina gets deposited into the empty spaces during passivation stage after ion-milling (alumina is deposited to passivate the sidewalls of MTJs.) This leads to thick alumina sidewalls, which can prevent solvents from reaching the lift-off layer. CO₂ snow cleaning process – utilizes expansion of either liquid or gaseous carbon dioxide through an orifice, which leads to the nucleation of small dry ice particles and a high velocity gas carrier stream – used to induce cracks on the sidewalls to facilitate lift-off process but, it fails to work for too thick alumina sidewalls. In addition, any residue e-beam resist and alumina on top of pillars can lead to poor or open electrical contacts.

On the other hand, if the etched depth is shallow, the RC constant can be large and may limit high-speed measurements.

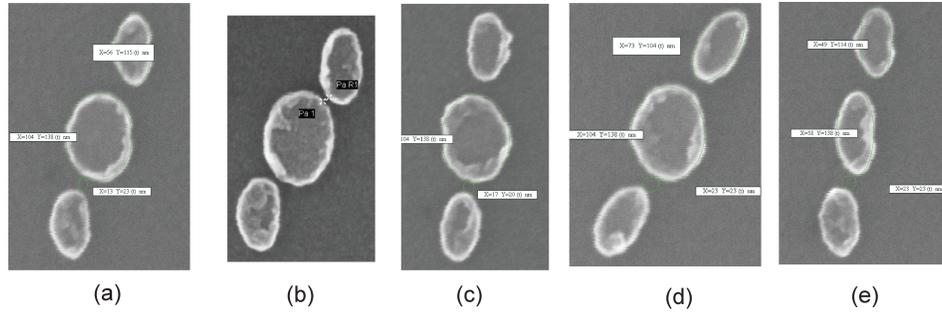


Figure 6.6. MCSTD gate designs for different logic functionalities. Different logic functions are derived from the difference in dipole couplings between the input and output device. Different input MTJ locations, angles and the output MTJ aspect ratios are tested.

The measured gap sizes (between the MTJs) after ion-milling were smaller than the nominal design spacing. This was due to the finite size of the electron beam and redeposition during the mask layer preparation process. Both effects increase the mask layer size and produce smaller spacing between the MTJs, which increases the magnetic coupling inside MCSTD gates. The CAD design for MCSTD e-beam lithography should take this effect into account and start from a larger gap size design to achieve the intended gap spacing.

6.7 Effect of shape irregularities

Since the devices interact via their fringing fields and utilize magnetic shape anisotropy, it is very important to maintain the uniformity in device shape and distance between the input and output MTJs. Figure 6.6 shows SEM images of various MCSTD gate designs. Although, some rough irregularities are observed in the line edges and shapes of MTJs, initial experimental results suggest that such processing limitations are well within the current technological limits of our processing. Furthermore, adopting other means such as perpendicular magnetic anisotropy instead of shape anisotropy to magnetically couple the devices can mitigate the device sensitivity to size and shape.

6.8 Electrical contacts

Finally, electrical contact lines were placed on the top of each MTJ device. Due to the proximity of the contact lines, this process was carried out by a combination of e-

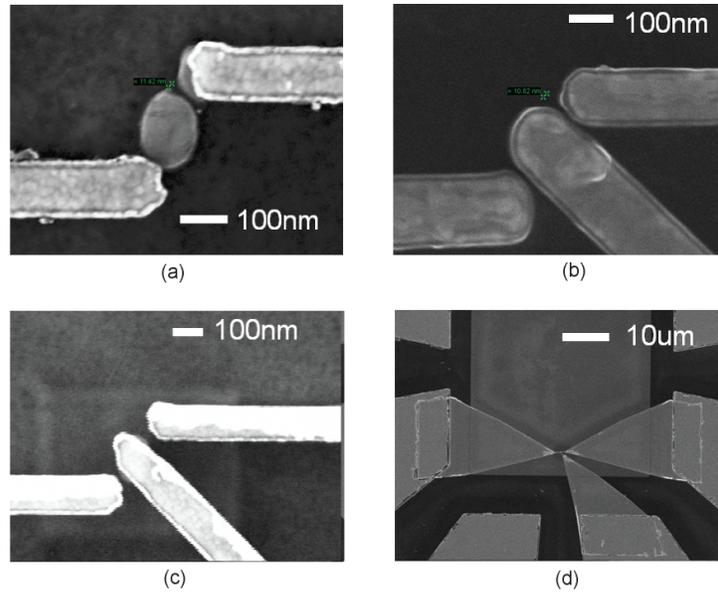


Figure 6.7. **SEM images of electrical contacts** on the output and input MTJs of MCSTD gates

beam and optical lithography (Fig. 6.7).

First, do ion-milling to clean the surface, i.e. etch away capping layer. Next, use e-beam lithography to pattern contact leads. No double layer is used for final lift-off process. Only PMMA is coated. After e-beam resist development, Au+Ta (Ta is adhesion layer) layers are deposited. When PMMA is lift-off, it may leave some “fencing” of Au+Ta layer.

Extended finger shape contacts were used for the input and output MTJs to avoid electrical shorts. Thanks to the small area of the electrical contacts, the overlap capacitance, which decides the RC delay of the measurements, could be minimized.

Reference

1. S. S. P. Parkin, "Spin-Polarized Current in Spin Valves and Magnetic Tunnel Junctions," *MRS Bulletin* 31, (2006)
2. A. A. G. Driskill-Smith, et al., "Electron-beam lithography for the magnetic recording industry: Fabrication of nanoscale (10nm) thin-film read heads," *Microelectronic Engineering*, 73-74 547-552 (2004)
3. Personal communication with Charles Rettner (2010)

Chapter 7.

**Experimental measurements of
Magnetically Coupled Spin-Torque
Devices (MCSTD)**

7.1 Introduction

In order to build a logic device, it is desirable to have a “gate” that changes the internal state of the device as a response to input signals. In MCSTD, the input spin-torque devices serve this role: whenever the input signal changes the magnetizations of the input MTJs, magnetic fringing fields that bias the output MTJ change. Then, the energy barrier height of the output MTJ – the internal state of MCSTD – is modulated. To validate this device mechanism, a series of measurements were conducted. Figure 7.1 illustrates the types of measurements and the parameters that are to be obtained from the experiments.

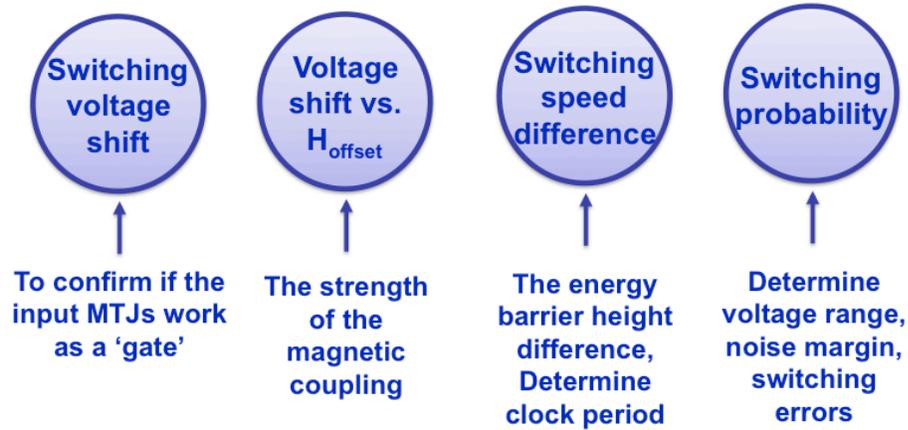


Figure 7.1. MCSTD measurement sequence and measurement parameters

7.2 Samples and Measurement setup

As introduced in Chapter 6, MCSTDs were fabricated on MTJ film stack that is deposited on Si substrates using a combination of ion-beam and magnetron sputtering at ambient temperature. The MTJ film structures used are (from bottom to top) replicated here for convenience.

100Å Ta | 300Å Ir₂₂Mn₇₈ | 6Å Co₄₀Fe₄₀B₂₀ | 30Å Co₇₀Fe₃₀ | 8Å Ru | 27Å Co₇₀Fe₃₀ | 8Å Mg | 4Å Mg in (95 Ar/5 O₂) | 20Å Co₄₀Fe₄₀B₂₀ | 50Å Ta | 50Å Ru

100Å Ta | 300Å Ir₂₂Mn₇₈ | 6Å Co₅₆Fe₂₄B₂₀ | 24Å Co₇₀Fe₃₀ | 5Å Ru | 27Å Co₇₀Fe₃₀ | 8Å Mg | 3Å Mg in (95 Ar/5 O₂) | 20Å Co₅₆Fe₂₄B₂₀ | 50Å Ta | 50Å Ru

Two film structures are different in CoFeB composition, MgO thickness, Ru thickness, CoFeB (free layer) thickness.

Figure 7.2 shows the measurement equipments that were used. The probe station that is shown in the middle is equipped with magnetic coils that can generate in-plane magnetic fields up to 2000 Oe. An AC pulse signal (pulse width 1~200ns) for high-speed measurement is generated with an Arbitrary Wave Generator. A sourcemeter is used to generate DC signal (long pulses in ms range). Although bias Tees are usually used to combine AC and DC signals for biasing purposes, we used it to remove high frequency electrostatic discharge (ESD) from DC signals. ESD is considered as one of reasons for MgO tunneling barrier breakdown. Other than using bias Tee, the first thing to do to prevent ESD is to properly ground yourself. Wearing a wristband with conductor usually does this. From my experience, wristband is not enough. Please make sure if your shoes are grounding yourself enough. Wearing shoes with thick rubber on the soles of shoes did not ground me sufficiently and I was blowing up dozens of MTJ tunneling barriers. Lastly, RF probes were used for high-speed measurements. To capture the switching

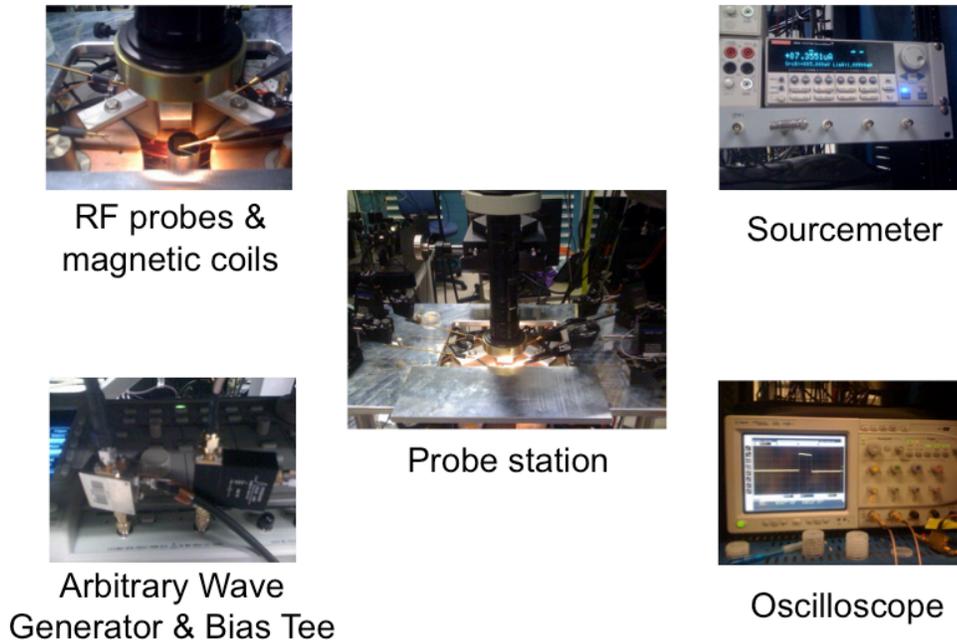


Figure 7.2. **Measurement setups** for MCSTD gate switching

voltage waveforms of the output MTJs, two RF probes were used to connect the device and an oscilloscope akin to a transmission line setup.

7.3 Input-dependent switching voltage shift

As mentioned in previous chapters, a MCSTD is a group of three spin-torque devices. We call the two smaller spin-torque control devices *input devices* and the larger center switching spin device *the output device*: the main idea is that the fringing fields from the input devices induce a change in the energy barrier height of the output device enabling it to predictably and reliably switch (when current is sent through the output MTJ). Fig. 7.3 explains this concept with measured device properties. As shown in Fig. 7.3(a), logical ‘1’ is defined to magnetize the input MTJ upward and logical ‘0’ in downward direction when its free layer is seen from the top. Fig. 7.3(b) illustrates the TMR versus magnetic field hysteresis plot of an output MTJ. As one sweeps magnetic fields, TMR changes between high and low resistance values along the directions of black arrows. This hysteresis loop is usually not centered around zero magnetic field because interlayer coupling (between the fixed and free layers) is not completely eliminated in actual devices. This interlayer coupling causes asymmetric switching magnetic fields or voltages as shown in Fig. 7.3(b) : $H_{\text{field}}(\text{AP} \rightarrow \text{P})$ is much smaller than $H_{\text{field}}(\text{P} \rightarrow \text{AP})$ in magnitude. Fringing fields from the input MTJ that MCSTD takes advantage of have the similar effect. They shift the TMR versus H_{field} loop depending on their directions, shifting switching voltage points. Since H_{fringing} from the input MTJs change their directions and magnitudes depending on input values, hysteresis loop shifts in different directions. For example, hysteresis loop shifts in the opposite directions for (spin up, spin up) and (spin down, spin down) inputs as shown in Fig. 7.3(c) and (d). This is why we should see the switching voltage shifts depending on the input signals. Fig. 7.3(e) plots experimentally measured dynamic resistance dV/dI versus voltage plots for different input signals. Switching voltage points are different for different inputs reflecting the difference in induced energy barrier heights. dV/dI values decrease as voltage increases due to Joule heating.

Figures 7.4 and 7.5 show repeated measurements (~40 iterations) of the output MTJ switchings for different input signals. They clearly demonstrate the input-dependent

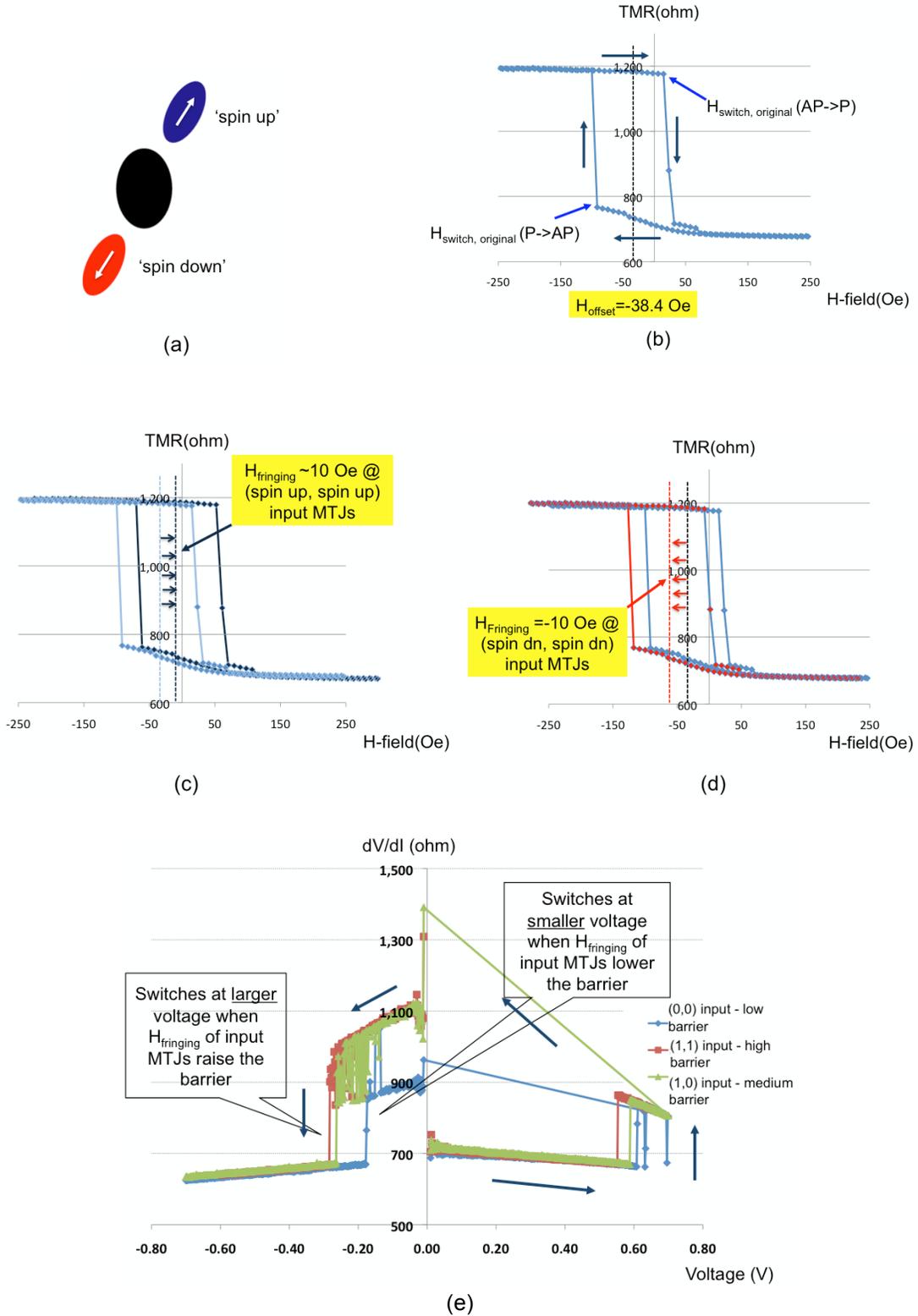


Figure 7.3. **Mechanism of switching voltage shift** in MCSTD (a) spin 'up' and 'down' definition in this discussion (b) single MTJ TMR vs. H-field hysteresis plot (c) the output MTJ TMR vs. H-field hysteresis curve at (spin up, spin up) inputs (d) the output MTJ TMR vs. H-field hysteresis curve at (spin down, spin down) inputs (e) the output MTJ dV/dI vs. Voltage hysteresis curve at various input signals

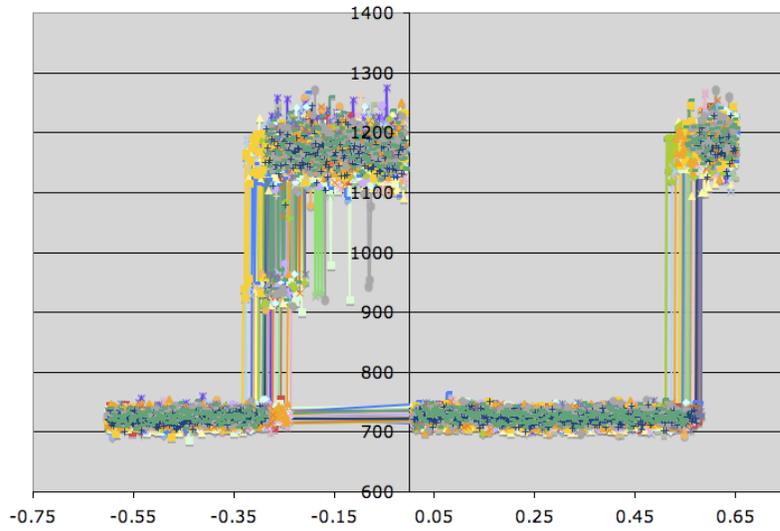
switching voltage shifts. TMR versus switching voltage hysteresis loops of MCSTD gates are illustrated. To avoid Joule heating effect, DC voltage pulse (100~300ms) is used to measure resistance. In Fig. 7.4, (a),(b) and (c) plots illustrate switching voltage points for (1,1), (1,0) and (0,0) inputs respectively. There are some deviations in switching voltage for the same input signal. Its standard deviation will be discussed in Sec. 7.4. In addition, magnetic reversal resulted in an intermediate value between high and low TMR values. This could originate from irregularities in fabricated device shape or magnetic properties that produced intermediate TMR state that traps the magnetizations before complete switching.

As shown in the plots, there is a clear shift in the hysteresis loop when the magnetizations of the input MTJs change. For the anti-parallel (AP) mode to parallel (P) mode switching, the average shift in the voltage was 100 mV and for the P->AP switching, it was 40 mV. In both cases, the voltage shifts are greater than the standard deviation of the switching voltage shift, which was measured to be 26.5mV. Switching voltage shifts were often asymmetric (between AP->P and P->AP) as shown in Fig. 7.4. A process variation induced device mismatch between the left and right input MTJ could be considered as a possible reason. In other devices shown in Fig. 7.5 demonstrated greater shifts in switching voltage points, such as 162 mV (Fig. 7.5(a)) and 162 mV (Fig. 7.5(b)).

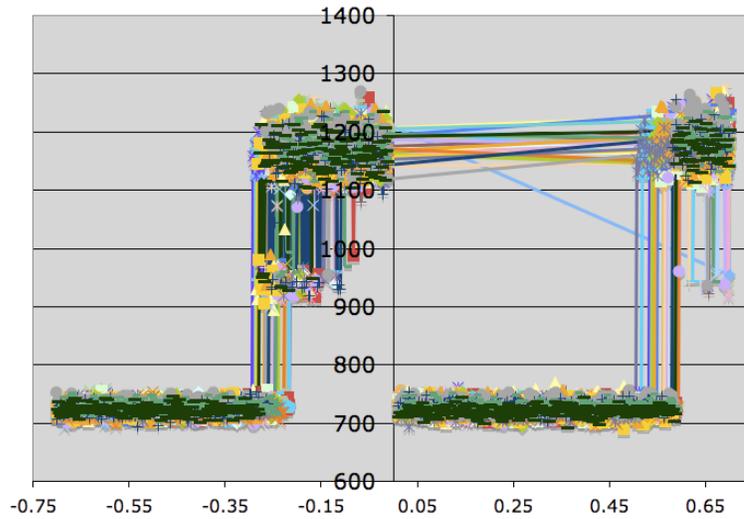
7.4 Input dependent switching voltage shift vs. H_{noise} and H_C

As a second step, we investigated the strength of magnetic coupling between the input and output MTJ. This experiment is to see if there is any dependence in the switching voltage shift on the external magnetic fields. If there were, the magnetic coupling between the input and output MTJs due to fringing fields are weak and can be easily affected by background thermal magnetic noise.

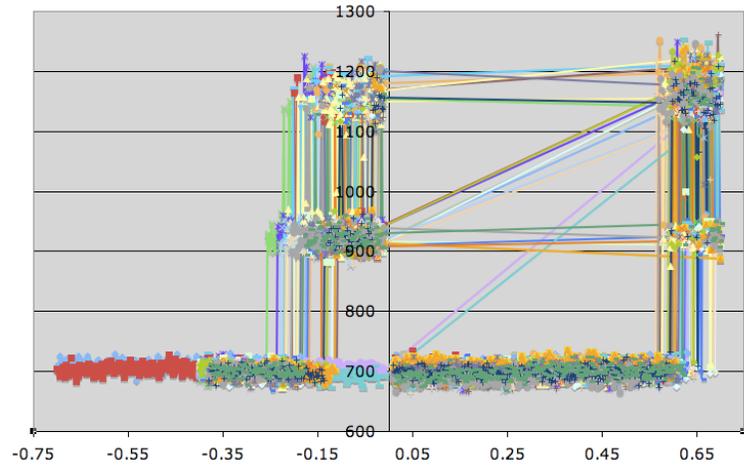
For the test on the magnetic coupling strength, the same switching voltage shift measurements as Sec. 7.3 were conducted at various background magnetic fields. Magnetic fields smaller than the coercivity, H_c , of the output MTJ, i.e., 10~70 Oe were applied as biasing magnetic fields. Figure 7.6, illustrates the results. There are two sets of parallel lines, each representing the switching from AP to P and vice versa. In both cases,



(a)

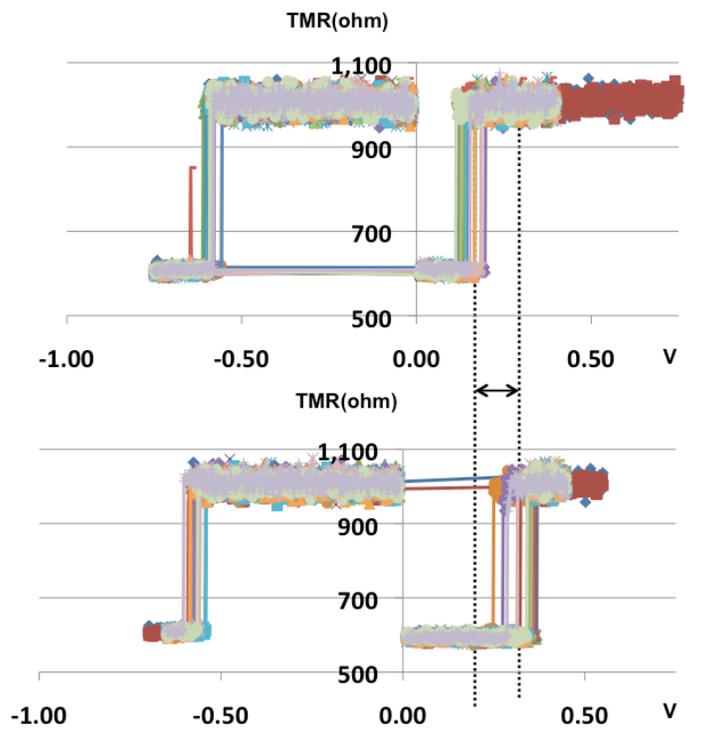


(b)

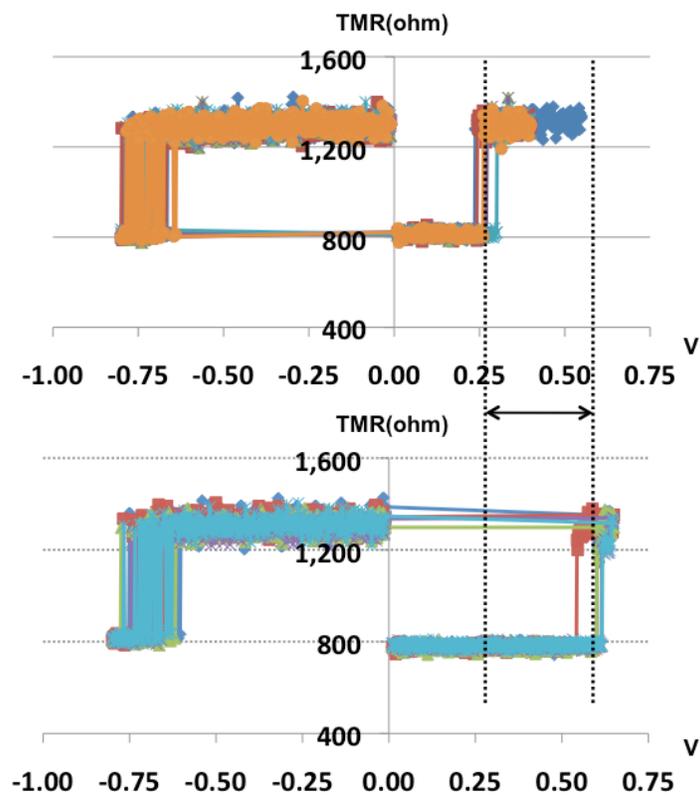


(c)

Figure 7.4. **Input signal dependence of switching voltage.** The switching voltage shifts when the input MTJ magnetization changes. (a) (spin up, spin up) = (1,1) input (b) (1,0) input (c) (0,0) input (device no. = s3256-12 d309)



(a)



(b)

Figure 7.5. **Input signal dependence of switching voltage.** The switching voltage shifts when the input MTJ magnetization changes. (a) device no. = s3256-12 d505 (b) device no. = s3256-12 d511. Upper plot illustrates (1,1) input and bottom plot (0,0) input

there was no dependence in the switching voltage shifts on the background magnetic fields. In other words, nearly constant switching voltage shifts discussed in Sec. 7.3 were observed throughout the noise magnetic fields considered. This result demonstrates that the magnetic coupling between the input and output MTJs is strong enough and the MCSTD device operation does not require any external magnetic field biasing.

This experiment provides us with useful information on the magnitude of the net fringing fields or magnetic coupling generated by the input MTJs. From the slope of the switching voltage in Fig. 7.6, one can deduce that it requires 10~15 Oe of external magnetic field to shift the switching voltage by 100 mV. Since the MCSTD gate in Fig. 7.6 shows a switching voltage shift of 80mV in AP-> P case, the fringing fields from the input MTJs is estimated to be 10 Oe in this case. Considering the fact that this particular output MTJ has a coercivity of 70 Oe, 10 Oe of fringing fields from the input MTJ amounts to 14%, which is strong enough to modulate the magnetic energy barrier.

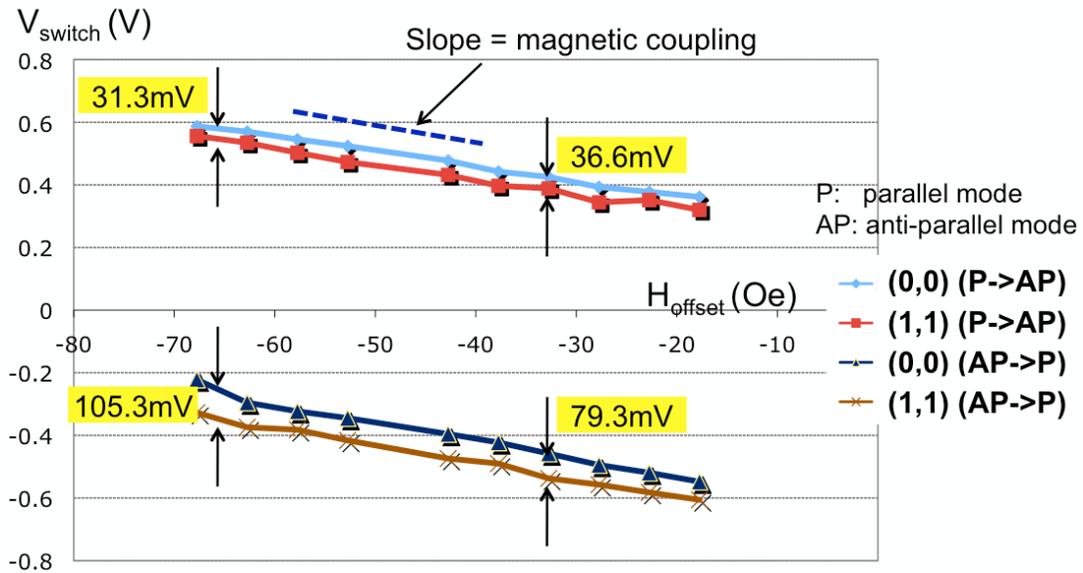


Figure 7.6. **Voltage shift of MCSTD gate versus H_{noise} .** No voltage shift dependence on background magnetic fields found. Magnetic coupling strength extracted from the slope. (device no. = s3256-12 d309)

In order to further understand the relationship between the magnetic coercivity and the successful MCSTD device operation, the switching voltage shifts versus the output MTJ coercivities are plotted in Fig. 7.7. We can induce from this plot that the coercivity

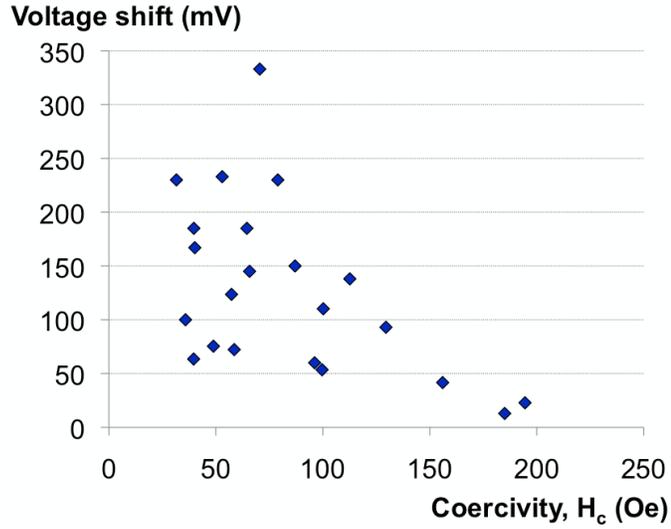


Figure 7.7. **Input dependent switching voltage shift vs. H_c .** Voltage shift decreases with H_c increase. H_c should be limited for successful MCSTD operation

of the output MTJ has to be maintained below a certain limit to maintain the input MTJ control on the switching characteristics of the MCSTD gate. In Fig. 7.7, the maximum switching voltage shift is inversely proportional to the output MTJ coercivity. The switching voltage shifts can be as large as 200 mV at $H_c < 100$ Oe but, it decreases as the H_c increases and becomes smaller than the average thermal fluctuations. This can be explained by theoretical J_c given by

$$J_c = \frac{1}{P} \frac{2e}{\hbar} \alpha M_s d (\pm H_{ext} + H_K + 2\pi M_s)$$

where P is spin polarization factor, α is the damping constant, M_s is the saturation magnetization, d is the thickness of free layer, H_{ext} is the external field, and H_K is the anisotropy field of the free layer [1]. J_c is linearly proportional to H_K , which is related to the coercivity. This linearity will make the fringing field from the input MTJs to be less obvious when the coercivity increases.

Magnetic coercivity depends on the material and microstructure properties. MTJs built with the same film stack and have the same structure should have the same

coercivity. Fig. 7.7 includes the data set from multiple wafers but, we still see some variations in the coercivity from the devices in the same wafer. A possible explanation is that the process variations such as line edge roughness (LER), etc. are inducing magnetic pinnings during the magnetic reversal process and this leads to an increase in coercivity. In addition, the fact that the switching voltage shifts in the coercivity < 100 Oe region are spread over 50 to 300 mV should also be attributed to the process variations in this prototype device. This gives us an idea how much fabrication defects or process variations we have and will be improved in the future.

7.5 Energy barrier height measurements & switching speed

In this section, the energy barrier height of the device is indirectly measured. We followed the approach that can be found in [2]. From this measurement (Fig. 7.8), the energy barrier height of the output MTJ was lowered from $42.89 k_B T$ to $30.84 k_B T$ due to the input MTJ magnetization that favors the output MTJ switch. The ratio between the high and low energy barrier height is 1.39:1.

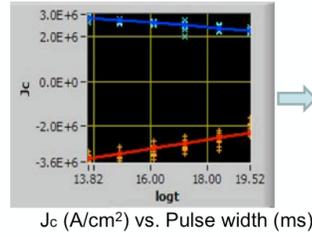
$$J_C = J_{C0} \left\{ 1 - (k_B T / E) \ln(\tau_p / \tau_o) \right\}$$

$$J_{C0} = \alpha \gamma e M_S t (H_{ext} \pm H_k \pm H_d) / \mu_B g$$

$$g = P / [2(1 + P^2 \cos \theta)]$$

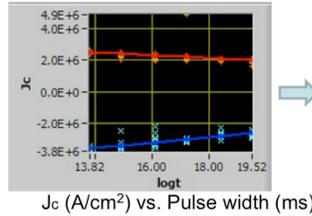
$$\left(\begin{array}{l} \alpha : \text{damping coeff.} \\ \gamma : \text{gyromagnetic coeff.} \\ t : \text{thickness of the free layer} \\ H_{ext} : \text{the external magnetic field} \\ M_S : \text{the saturation magnetization} \\ P : \text{the spin polarization} \\ H_d : \text{the out - of - plane magnetic anisotropy} \\ \theta : \text{the angle between the free/fixed layer} \end{array} \right)$$

(1,1) input case:



Energy barrier, $E/k_B T$:
 AP->P : **30.84**,
 P->AP: **44.2**

(0,0) input case:



Energy barrier, $E/k_B T$:
 AP->P : **35.73**,
 P->AP: **42.89**

Figure 7.8. **Energy barrier height of MCSTD gates** extracted from J_c vs. pulse width relations

Energy barrier height modulation leads to a difference in switching speed. In other words, the same device demonstrates different switching speed if the input MTJ magnetizations change. Fig. 7.9 shows the switching speed of MCSTD for (1,1) input – low barrier and (0,0) input – high barrier cases. Switching speed of the two cases can be considerably different at low voltages. For example at voltage 0.95V, switching time is $T_{high\ barrier} : T_{low\ barrier} = 75ns : 20ns = 3.75:1$. By choosing a circuit clock period that is shorter than $T_{high\ barrier}$ but longer than $T_{low\ barrier}$, for example, clock period = 50ns, we can virtually prevent the switching of MCSTD with high barrier case. One of the critical disadvantages of spin-torque devices such as MTJs as a logic device is that I_{on}/I_{off} or TMR is less than 5. However, if we use the switching speed difference and judiciously adjust the clock period that the spin-torque devices are driven, we can make the device in the high barrier case virtually impossible to switch.

The switching speed difference diminishes as the voltage becomes high, i.e., the switching energy is large enough to overcome the high barrier as well. This issue can be addressed by enhancing the magnetic coupling between the input and output MTJ such that the switching speed difference becomes much greater. The shortest switching time of MCSTD in our experiment was 6ns as shown in Fig. 7.9. However, it was mainly due to

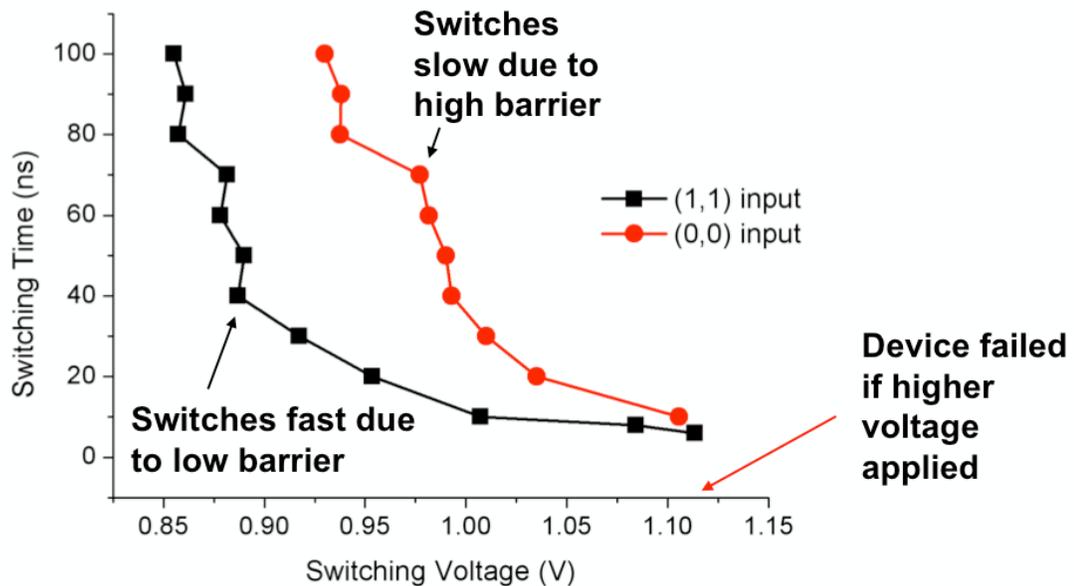


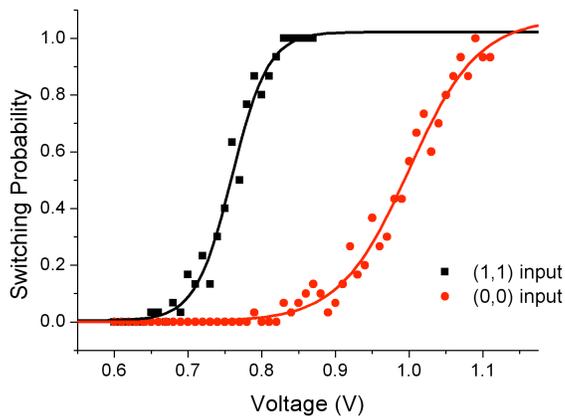
Figure 7.9. Input signal dependent switching speed difference

the resolution limit in measurement setup and reliability issues in MTJ free layer. With higher resolution equipment and emphasis on the spin-torque device reliability, faster switching speed will be immediately achieved.

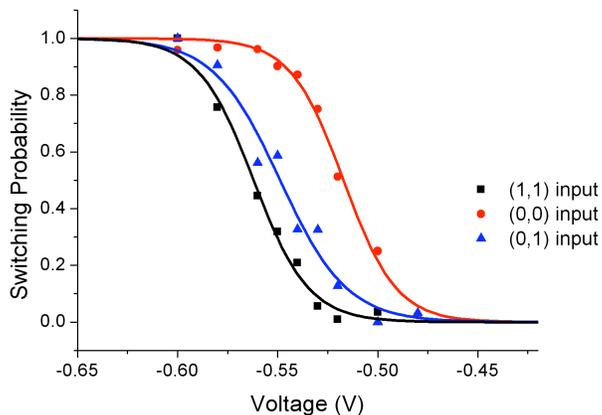
7.6 Switching probability and logic operations

Figure 7.10 illustrates the switching probability plots for three different MCSTD gates. All devices demonstrate “input-signal dependent” switching probabilities. In other words, the same device shows different switching characteristics depending on its input signals. This is another strong evidence that the input MTJ magnetizations and induced fringing fields successfully work as a “gate” and energy barrier height modulation. As the switching voltage is increased, the device with lower energy barrier height will switch earlier than the one with higher barrier height. Fig. 7.10 (b),(c) show that (0,1) and (1,0) input cases all come in between the (0,0) and (1,1) input cases. This result reflects that their energy barrier heights are in between those of (0,0) and (1,1) input cases as expected.

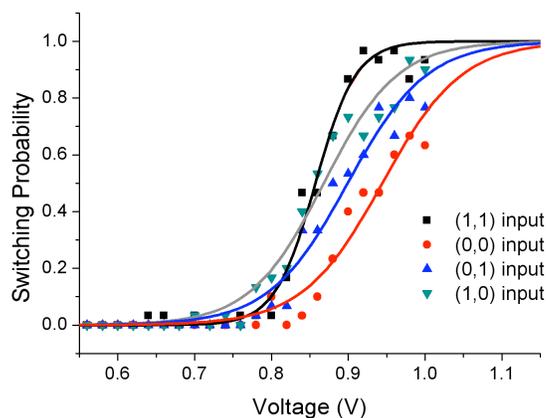
One of the important design advantages of MCSTD based non-volatile logic is the ability to control the switching probability and hence logic error rate. If the MCSTD logic



(a)



(b)



(c)

Figure 7.10. **Experimental demonstration of switching probability change versus the input device magnetization directions.** MCSTD gate logic operation is based on the change in its switching characteristics depending on the input device magnetizations. (a) Large gap between the low barrier ((1,1) input) and the high barrier ((0,0) input) guarantees low error rate in computations. (b),(c) (0,1) and (1,0) input cases all come in between the (0,0) and (1,1) input cases as expected

gate shown in Fig. 7.10 (a) is biased at 0.84V, the gate switches at 100% accuracy at (1,1) input and does not switch at (0,0) with 97% accuracy, which makes the overall error rate less than 3%. This is remarkably lower error rate compared to other nanomagnet logic device works, such as Magnetic Quantum Cellular Automata [3], whose minimum error rate is reported to be 15% in even the best case simulations [4]. The error rate of MCSTD can be made even smaller by 1) separating the switching probability plots farther apart and 2) making the slope of the curve steeper. The device in Fig. 7.10(c) needs to improve more on these aspects because the separation of the switching probability plots are narrow and some cross over in certain region. Each task can be achieved by increasing the energy barrier height ratio between the high/low barriers and better confining the magnetic flux between the input and output MTJs.

Switching probability versus voltage plot can be used to identify the logic functionality that can be performed by a MCSTD gate. For example, as shown in the truth table for NAND device (Table 7.1), switching should happen only at (1,1) input and not at the others. This requires the switching probability plots to look similar to Fig. 7.11

	Input 1	Input 2	Output	Need switching?
NAND	0	0	1	No
	0	1	1	No
	1	0	1	No
	1	1	0	Yes
NOR	0	0	1	No
	0	1	0	Yes
	1	0	0	Yes
	1	1	0	Yes

Table 7.1 Truth table for NAND and NOR logic. Assumed that Output=1 at initial state

upper plots: switching probability plots for (0,1) and (1,0) inputs should be placed near (0,0) input case and keep some distance from (1,1) inputs. By arranging the operating voltage range that is in MCSTD gate will switch only for (1,1) input and function as NAND gate. Figure 7.12 shows the experimentally measured time response of NAND MCSTD gate. This device does show different switching voltage points for different inputs: from 0.28V to 0.40V, switching only happens for (1,1) input and not for other inputs, which is a NAND gate. In contrast, for a NOR MCSTD gate, the switching probability should resemble that of Fig. 7.11 bottom plot: the switching probability plots for (0,1) and (1,0) inputs should be placed near that of (1,1) input case. Hence, the biasing voltage range inside the red dotted box is chosen, MCSTD gate will switch for (1,1), (0,1) and (1,0) inputs and not for (0,0) input. This type of time response switching behavior is shown in Fig. 7.13. This MCSTD gate switches for (1,1), (0,1) and (1,0) inputs and not for (0,0) input at the voltage range of 0.52~0.58V. MCSTD with (0,0) inputs also started to switch at the biasing voltage greater than 0.58V. As shown in these two examples, one of the great advantages of MCSTD is that different logic device can be made by simply changing its shape, the location and angle of the input device, which will lead to low manufacturing cost.

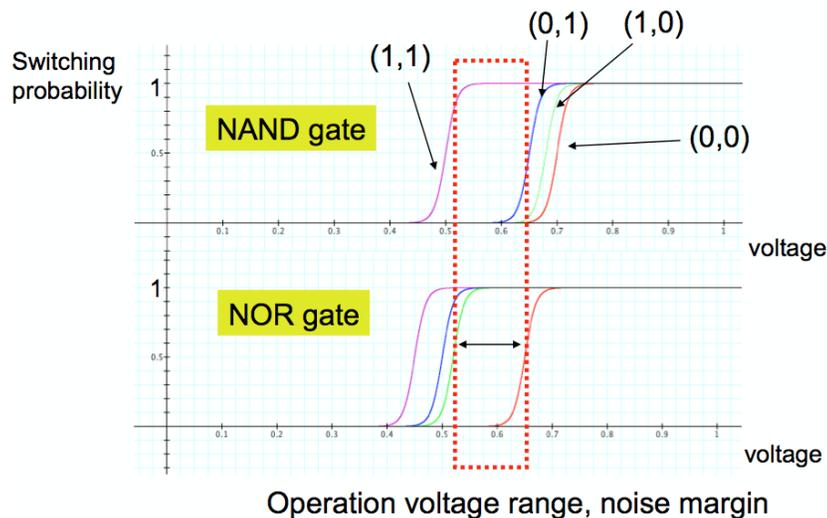


Figure 7.11. **Proposed voltage operation range of MCSTD gates.** NAND and NOR gates have different switching voltage for different inputs, which makes them function as different logic devices. The numbers in parentheses are input signals.

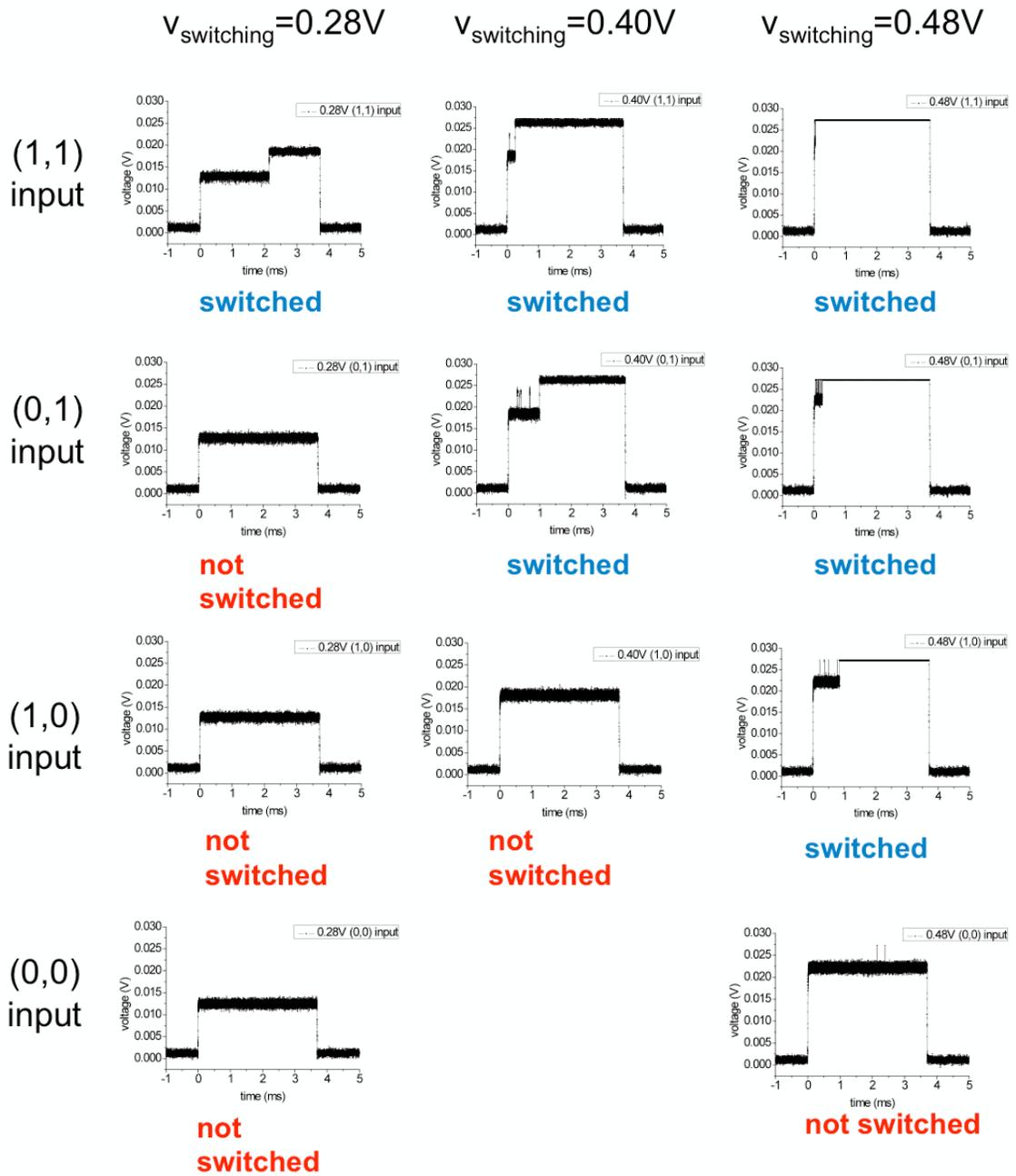


Figure 7.12. Experimentally measured time response of NAND MCSTD gate. Each row represents four different input signals. Same biasing voltage is used for the same column. Operates as NAND gate at 0.28V~0.40V

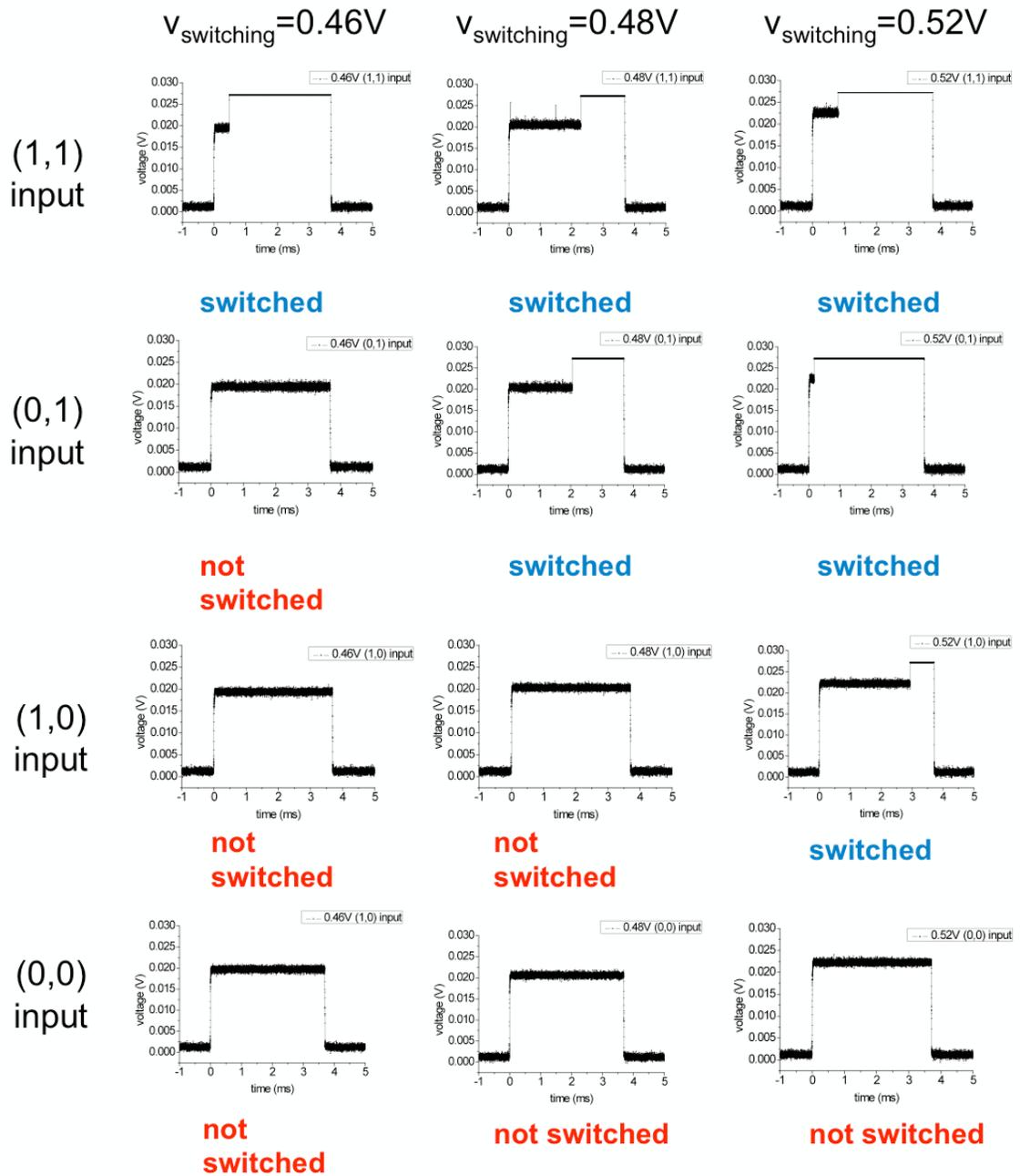


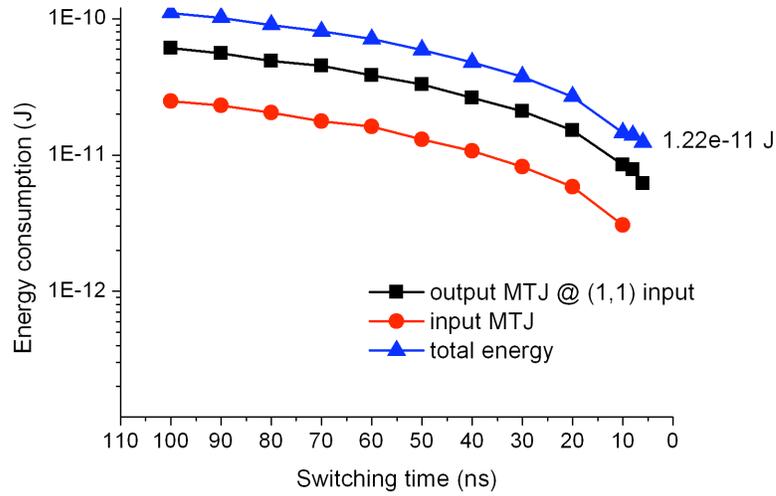
Figure 7.13. Experimentally measured time response of NOR MCSTD gate. Each row represents four different input signals. Same biasing voltage is used for the same column. Operates as NOR gate at 0.52V~0.58V

7.7 Energy consumption measurements

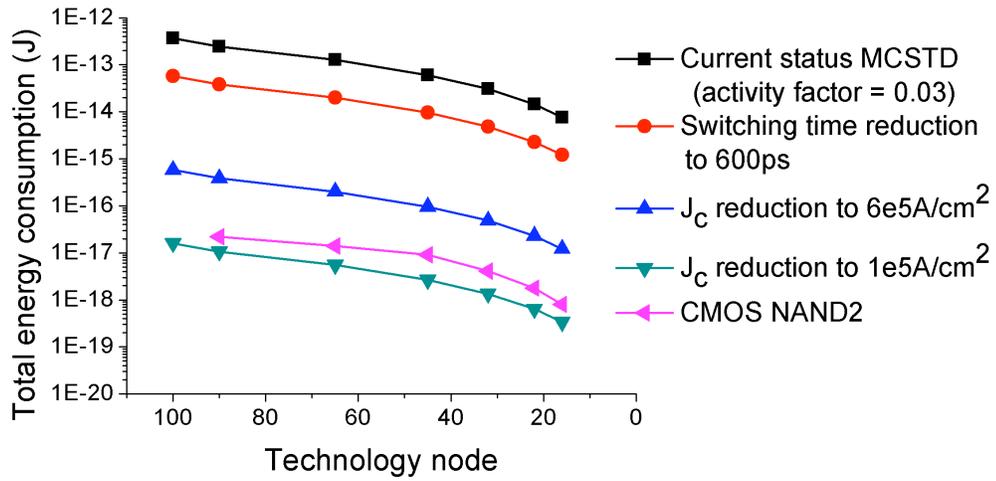
Using TMR versus voltage plots and switching time data, the energy consumption of MCSTD gates are measured. Since MCSTD is basically a group of MTJs, the energy consumption is also the sum of energy consumption in individual MTJs. The total energy consumption calculation should include energy consumption from interconnection and CMOS peripherals and it was illustrated in Sec. 5.1. In this section, we consider energy consumption from MCSTD gate only. To make this a fair comparison, we considered the energy consumption from CMOS gates only as a reference.

From data points in Fig. 7.10, we calculate the energy consumption of the input and output MTJs. From the various switching time and switching voltage, corresponding consumed energies are calculated. Lowest energy consumption point of 1.22×10^{-11} J corresponds to shortest switching time of 6ns. In order to project the achievable energy consumption from this data point, a number of extrapolations are made. First, device dimension is scaled down to 15nm. The input MTJ can be considered as a scaled down version of the output MTJ and suggests how much energy scaling we can get from device scaling. In addition, it was assumed that the switching time is reduced to 600ps, which is an order of magnitude shorter than measured value. Still MCSTD gate consumes three orders of magnitude higher than that of CMOS NAND gate. MCSTD becomes more energy efficient than CMOS only when the current density scales down to 1×10^5 A/cm². Spin-torque transfer based technology makes one of the most energy efficient memory device but, for logic, it is still less competitive than the highly scaled CMOS gate. Reduction in switching current density, spin-torque device resistance and switching time are required.

Ultimate scaling limit of ferromagnetic device is the super-paramagnetic limit, where the ferromagnetic devices lose hysteresis. As shown in the following equations, magnetic devices are designed to have data retention time over 5~10 years, which imposes the device scaling limit. However, 5~10 years of data retention time is mainly tailored for memory device. Logic device would not require these long time of data retention time. Perhaps several weeks or a month will be enough for logic application. This trade-off between thermal stability with device scaling could be a possible approach to further reduce energy consumption in MCSTD gates.



(a)



(b)

Fig. 7.14. (a) Total energy consumption of MCSTD gate calculated from experimental data in Fig. 7.9. (b) Energy consumption projection of MCSTD gate

$$\frac{1}{\tau} = f_0 e^{-\frac{\Delta E}{k_B T}}$$

$$\Delta E = K_u V$$

$$\frac{\Delta E}{k_B T} > 25$$

$$\Rightarrow 7.5 \text{ nm}(\text{width}) \times 9.5 \text{ nm}(\text{length}) \times 2 \text{ nm}(\text{thickness})$$

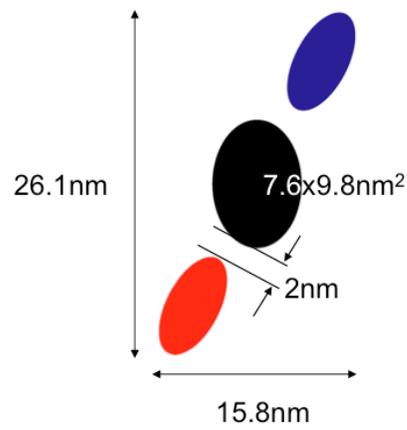


Fig. 7.15 Device dimension of MCSTD gate at superparamagnetic limit

Reference

1. W. Kim *et al.*, “Enhanced switching current density due to resonant precession in current-induced magnetization switching,” *Applied Physics Letters*, vol. 90, issue 212504 (2007)
2. Hayakawa, *et al.*, “Current-Induced Magnetization Switching in MgO Barrier Magnetic Tunnel Junctions with CoFeB-Based Synthetic Ferrimagnetic Free Layers,” *IEEE Trans. Magnetism*, vol. 44 (2008)
3. D.B. Carlton, N.C. Emley, E. Tuchfeld, J. Bokor, “Simulation of nanomagnetbased logic architecture,” *Nano Lett.* vol. 8, pp. 4173–4178 (2008)
4. F.M. Spedalieri, A.P. Jacob, D. Nikonov, V.P. Roychowdhury, “Performance of magnetic quantum cellular automata and limitation due to thermal noise,” arXiv:0906.5172v1 (2009)

Chapter 8.

Multi-scale Simulations of Partially Unzipped CNT Hetero-junction Tunneling Field Effect Transistor

A version of this chapter has been accepted and will be published as L. Leem, A. Srivastava, S. Li, G. Iannaccone, J. S. Harris, G. Fiori, “Multi-scale Simulations of Partially Unzipped CNT Hetero-junction Tunneling Field Effect Transistor,” *Proc. International Electron Device Meeting*, 2010

8.1 Motivation

In this chapter, we switch gears to change the fundamental design of the MOSFET and lower the sub-threshold slope (SS) below 60 mV/dec and continue to scale V_{th} and V_{dd} down. Band-to-band Tunneling Field Effect Transistors have recently attracted deep interest in the research community because of their small SS allowing reduced supply voltages for digital logic [1]. While a number of band-to-band tunneling transistors have been reported to achieve $SS < 60\text{mV/dec}$ [2], their potential applications are limited by 1) small I_{on} current as compared to conventional MOSFETs 2) ambipolar I-V characteristics and 3) $SS < 60\text{mV}$ is obtained only in a very limited V_{gs} interval. In this work, we present for the first time, multi-scale approach consisting of Density-Functional Theory (DFT), atomistic Tight-Binding (TB) and Molecular Dynamics (MD) calculations to simulate Carbon-nanostructures-based TFETs with type-II tunneling barriers that results from partially unzipping carbon nano-tubes (CNTs) [3]. The concept of partially unzipping CNT (Fig. 8.1) is promising because one can obtain 1) high quality GNRs with reduced edge roughness [3] and 2) type-II heterojunctions due to the difference in energy bandgaps and electron affinities between the primitive and unzipped CNT (Table 1). Band asymmetries induced by the heterojunctions can overcome the limitations of tunneling FETs (TFETs) by eliminating the I-V ambipolarity and further improving SS. Four different device configurations (Fig. 8.2) have been investigated for optimal device performance according to ITRS requirements [4].

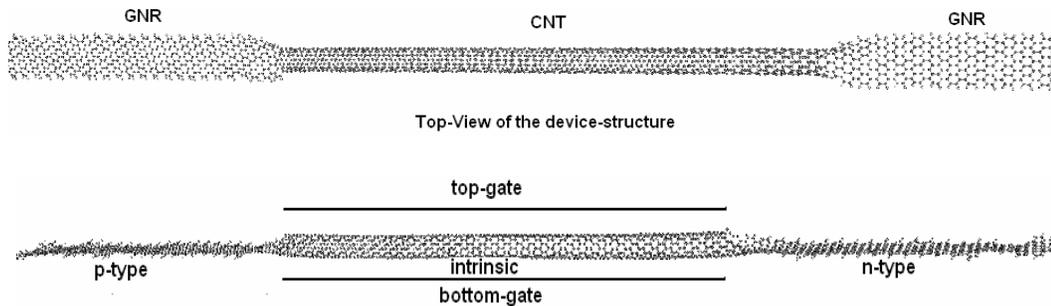


Figure 8.1. Geometrically relaxed partially unzipped Carbon nanotube (CNT) Molecular dynamic simulation was conducted to study the energy-relaxed configurations of GNR/CNT heterostructures after unzipping (2507 atoms, total length 37nm)

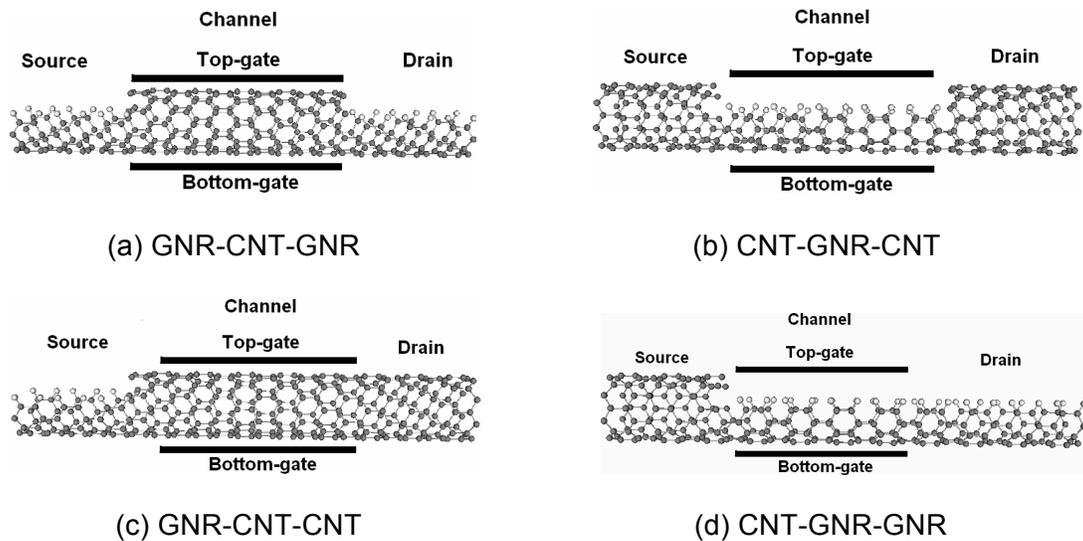


Figure 8.2. Cross-sectional schematic of simulated GNR/CNT heterostructures to study the effect of GNR and CNTs on the tunneling FET performance

8.2 Device simulation process

The flow diagram of our multi-scale approach is shown in Fig.8.3. First, energy bandgaps and electron affinities of GNRs and CNTs were computed using DFT, the bandgap-underestimation in DFT was corrected through Extended Huckel's Theory (EHT). Second, Molecular Dynamics simulations were performed to estimate the geometric

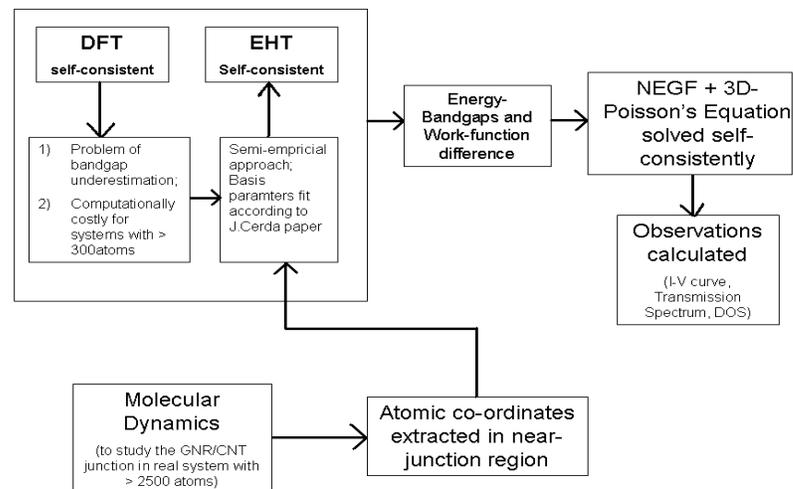


Figure 8.3. Flow chart of multi-scale simulations for heterojunction Tunneling FETs

CNT (N,0)	GNR (n)	E_g (CNT) (eV)	E_f (CNT) (eV)	E_g (GNR) (eV)	E_f (GNR) (eV)	$\Delta\phi$ (eV)	Type-II	Broken gap	$E_{g,eff}$ (eV)
14	24	0.64	3.7831	0.352	3.628	-0.0111	✓		0.3409
13	22	0.77	3.8215	0.46	3.618	-0.0485	✓		0.4115
11	16	0.79	3.745	0.624	3.582	-0.08	✓		0.544
10	16	1.01	3.813	0.624	3.582	-0.038	✓		0.586
8	16	1.054	4.017	0.624	3.582	-0.22	✓		0.404
7	14	1.184	4.201	0.047	3.567	-0.0655		✓	-0.0185

Table 8.1. Electronic parameters calculated using DFT and EHT theories

structure of the hydrogen-passivated partially unzipped CNT device structure. Finally, the 3D-Poisson equation within the non-equilibrium Green's function (NEGF) formalism and tight-binding methods were used to calculate I-V characteristics of the device.

8.3 Device operating principle and simulation results

The ambipolarity of TFETs occurs for symmetric energy bands (Fig. 8.4), where tunneling happens for both $V_{gs} > 0$ and $V_{gs} < 0$. Popular approach to suppress the ambipolarity is to break the symmetry with asymmetric doping or gate electrode location [5]. However, the effectiveness of these techniques is limited to narrow V_{gs} range [5]. Type-II heterojunctions can fundamentally break this symmetry, leading to excellent sub-

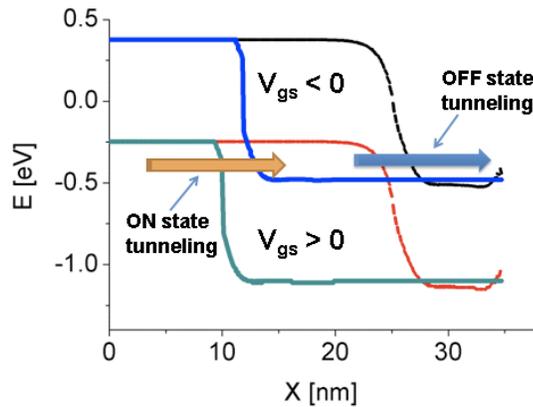


Figure 8.4. Symmetric energy bands in Homojunction TFETs
Tunneling happens for both $V_{gs} > 0$ and $V_{gs} < 0$

threshold region characteristics. Large bandgap material is used for channel and drain region to suppress tunneling in the device's off-state. At the same time, small bandgap material is used for source region to increase I_{on} . Type-II heterojunction is more effective than type-I due to smaller effective bandgap and larger energy barrier height at heterojunctions. For example, in GNR(*source*)-CNT(*channel*)-CNT(*drain*) Tunneling FET, tunneling is enhanced in the ON state due to smaller effective bandgap (Fig. 8.5(e)). In the OFF state, tunneling is strongly suppressed by larger bandgap channel material and high energy barrier of type-II heterojunction (Fig. 8.5(f)). Compared to the CNT and GNR homojunction TFETs, GNR/CNT heterojunction TFETs demonstrate superior sub-threshold region characteristics - $10^4 \times$ smaller I_{off} , 61% smaller $SS=22\sim 26\text{mV/dec}$ (compared with CNT-TFET) and $2.88\times$ wider voltage region, where ambipolarity is strongly suppressed (compared with GNR TFET) (Fig. 8.6). For fair comparisons, CNT and GNR TFETs were optimized for the best performance (Fig.8.7). Among four heterojunction configurations, GNR-CNT-GNR and GNR-CNT-CNT demonstrate better subthreshold slope and I_{off} than the rest due to large energy bandgap material in the channel (CNT(8,0)). If GNR $n=16$ is used for channel material, weak I-V ambipolarity was observed (Fig.8.8). Single (GNR-CNT-CNT) and double heterojunction configurations (GNR-CNT-GNR) are compared in Fig.8.10.

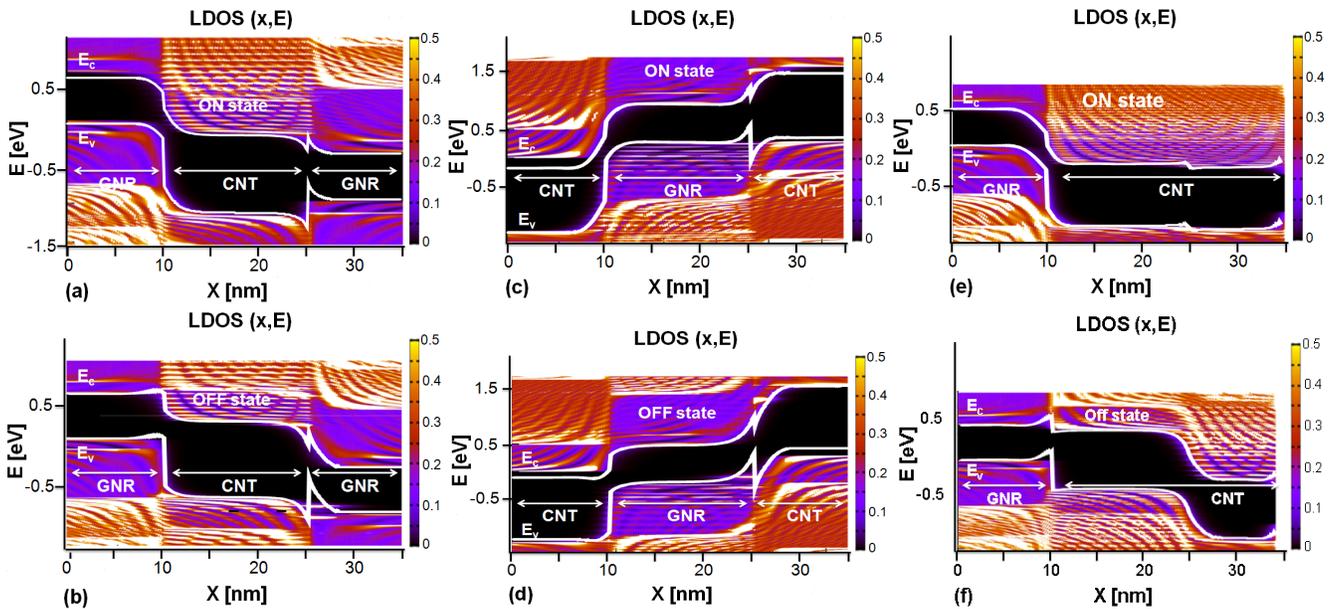


Figure 8.5. Band diagram (white solid lines) and local density of states of GNR/CNT tunneling FETs. GNR-CNT-GNR (left column), CNT-GNR-CNT (center) and GNR-CNT-CNT (right) configurations shown for ON (top row) and OFF state (bottom row). CNT(8,0) is partially unzipped to create GNR $n=16$ at source/drain or channel regions. Due to heterojunctions, symmetric energy bands present in homojunction TFETs are removed

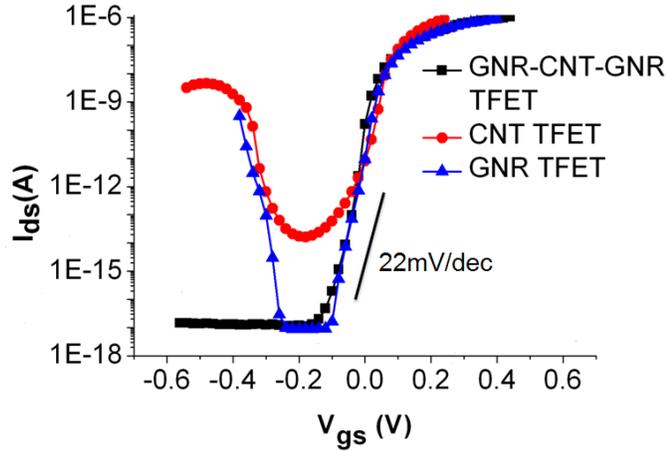


Figure 8.6. Comparison of I-V curve of Carbon based tunneling FETs. The workfunction of the gate electrode is adjusted that $I_{ds}=10^{-11}$ A is at $V_{gs}=0$ V. Compared to CNT TFET, GNR/CNT heterojunction TFET shows better I_{off} and smaller SS, while they show comparable performance with GNR TFET, with I-V ambipolarity nearly completely suppressed. This is because of the larger E_g of the material (CNT) in the channel

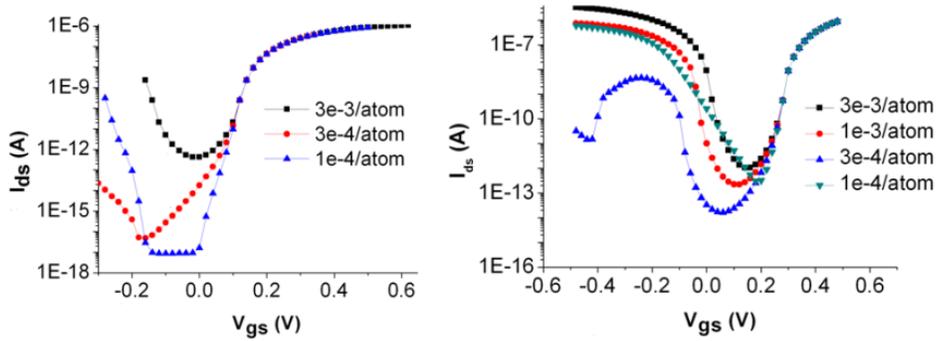


Figure 8.7. Ambipolar I-V characteristics of homojunction TFETs. GNR TFET (left) and CNT TFET (right) optimized for fair comparison with GNR/CNT Heterojunction TFET. Doping levels (in molar fractions) in drain region were reduced (as compared to source region) to suppress the ambipolarity. CNT(14,0) ($E_g=0.64$ eV) was chosen for CNT TFET to match the $E_g=0.624$ eV of GNR $n=16$

Subthreshold characteristics of GNR/CNT heterojunction TFET depends on the channel length, doping and V_{ds} (Fig. 8.11). Direct tunneling between source and drain dramatically increases with shorter channel length and as the V_{ds} and doping increase, the tunneling through channel-drain junction in the OFF state becomes stronger. I_{on} strongly depends on the injecting states from the source region. GNR-CNT-(GNR or CNT)

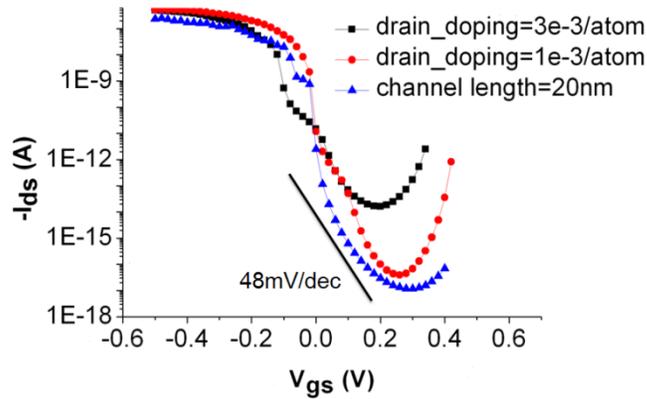


Figure 8.8. CNT-GNR-CNT heterojunction TFET I-V characteristics. When unzipped region (GNR) is used as a channel region, device turns on at $V_{gs} < 0$, effectively working as a PMOS device. However, due to smaller bandgap of GNR in the channel region, weak ambipolarity is witnessed

configurations show comparable I_{on} as those of GNR homojunction TFETs, whereas, CNT-GNR-(CNT or GNR) show similar I_{on} to CNT TFETs. However, this is a conservative estimation, since it can be further improved by an inherent stress developed in the junction region due to the partial unzipping of the CNT.

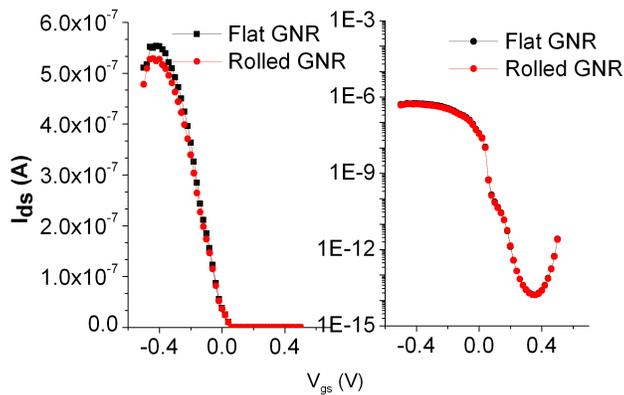


Figure 8.9. I-V comparison between completely unrolled GNR ("flat") and rolled GNR. I-V curves of the two cases are compared (left plot in linear scale and the right plot in log scale). Average difference in the I_{ds} is less than 4.6%, i.e., no "curvature" effect on I-V characteristics is found in partially unzipped CNT

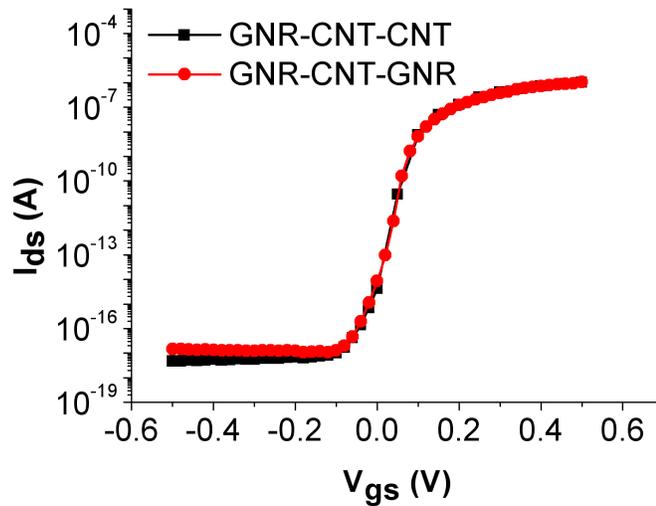


Figure 8.10. Comparison between single (GNR-CNT-CNT) and double heterojunction (GNR-CNT-GNR) TFETs. From I-V curve point of view, no distinct difference was found in I_{on} , I_{off} and subthreshold swing of these devices, indicating it is sufficient to have unzipping of CNT at source/channel interface only for this particular device dimension and operation voltage range

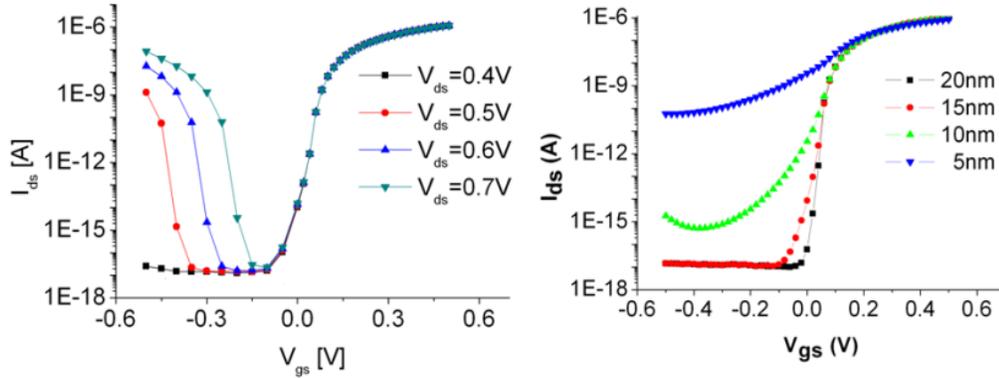


Figure 8.11. V_{ds} (left) and channel length (right) dependence of GNR/CNT Heterojunction TFETs. Subthreshold slope (SS) and I_{off} show dependency on the V_{ds} (left) and the channel length (right). Large V_{ds} increases the device off-state tunneling through the channel-drain junction. This tunneling disappears when $V_{ds} \leq 0.4V$ and removes the ambipolarity. Short channel length $< 15nm$ increases the direct tunneling between source and drain and brings about large I_{off}

	GNR/CNT heterojunction TFET	GNR	CNT
Contact resistance	Low. Width and chirality of GNR can be controlled by etching time, so that the workfunction can be modified to reduce the contact resistance	Low	High
Channel mobility	High. Using CNT as channel region can increase the mobility of carriers. Graphene and GNR are more sensitive to the surface ad-atoms than CNTs.	Low	High

Table 8.2. Fabrication benefits of GNR/CNT Heterojunction TFET vs. GNR, CNT TFET

8.4 Conclusion

We have investigated the performance of partially unzipped CNT heterojunctions, by means of a multi-scaled approach based on DFT, EHT, Molecular Dynamics (MD) and self-consistent tight-binding simulations of carrier transport. GNR/CNT heterojunctions demonstrated to be good candidates for low voltage logic applications and show better performance in terms of low subthreshold slope and strongly suppressed ambipolar behavior as compared to CNT and GNR TFETs.

References

1. J. Knoch, J. Appenzeller, “Tunneling phenomena in carbon nanotube field-effect transistors,” *Phys. Stat. Sol. (A)*. 205, no. 4, (2008)
2. W. Choi, B. Park, J. Lee, T.-J. K. Liu, “Tunneling Field-Effect Transistors (TFETs) With Subthreshold Swing (SS) Less Than 60 mV/dec,” *IEEE Electron Device Lett.* 28 (2007)
3. L. Jiao, L. Zhang, X. Wang, G. Diankov, H. Dai, “Narrow graphene nanoribbons from carbon nanotubes,” *Nature* 458, 877 (2009)
4. ITRS <http://public.itrs.net>
5. P. Zhao, J. Chauhan, J. Guo, “Computational Study of Tunneling Transistor Based on Graphene Nanoribbon,” *Nano Letter* 9, no. 2. (2009)

Chapter 9.

ERSA: Error Resilient System Architecture for Probabilistic Applications

9.1 Introduction

Reliability is a major concern for power-constrained computing systems in advanced CMOS technologies. Several mechanisms threaten to significantly increase the number of errors experienced by future systems – erratic bit errors, transient (soft) errors, early-life failures (infant mortality), transistor aging, Vccmin challenges, process variations and variations in operating conditions (e.g., voltage droops) [Agostinelli 05, Borkar 04, Van Horn 05]. To overcome these challenges, two classes of techniques have been employed in many systems:

1. Conservative design to ensure correct operation, e.g., guardbanding, conservative voltage scaling. Many of these techniques appear to be running out of steam [Gelsinger 06].
2. Fault-tolerant systems that detect and recover from errors through expensive redundancy.

The *Error Resilient System Architecture (ERSA)*, presented in this paper, leverages two emerging trends in future computing platforms:

1. Proliferation of multi- and many-core systems.
2. New killer applications, e.g., data mining, market analysis, cognitive systems and computational biology, which are expected to drive demands for computation capacity. Such applications are also referred to as Recognition, Mining, Synthesis or *RMS* applications [Dubey 05].

Unique properties of RMS applications include:

1. **Massive parallelism:** Massive amounts of data are processed to build mathematical models and to apply models to help answer real-world questions.
2. **Algorithmic resilience:** Unlike conventional computing, RMS applications are tolerant to imprecision and approximation to make analysis of complex systems tractable. In addition, an iterative approach is often used to refine computation results such that the effects of inaccuracies can be reduced by subsequent iterations.
3. **Cognitive resilience:** Computation results need not always be correct as long as the accuracy of the computation is “acceptable” to human users [Breuer 05].

Several aspects of algorithmic and cognitive resilience of various applications to low-order data bit errors have been addressed previously by several researchers [Breuer 05, Chakrapani 06, Hayes 07, Li 06, Shanbhag 02, Yu 00, Wong 06]. While such data error resilience is clearly a win, it alone is not sufficient for probabilistic applications to converge and generate useful results when executed on unreliable computing hardware. At high error rates, the errors in high-order bits of data and application control flow significantly affect application performance (details in Sec.2).

ERSA uses the following techniques to overcome the above challenges of high-order bit errors and control errors:

1. *Asymmetric reliability*, i.e., mixing processor cores of various “reliability levels” in many-core architectures.
2. Software optimizations including minimally-intrusive yet effective modifications to RMS algorithms.
3. Light-weight checks such as timeouts and memory bounds violation checks. ERSA does not rely on expensive error detection techniques.

Error injections in actual ERSA hardware platforms demonstrate that, even at extremely high error injection rates of 20,000 errors/sec/core or 2×10^{-4} error/cycle/core into architecturally-visible registers, ERSA delivers RMS applications with the following characteristics:

1. No system crashes.
2. 90% or better accuracy of output results (within cognitive resilience limits as demonstrated using actual applications).
3. Minimal execution time increase (20% or less).

While ERSA is optimized for probabilistic applications, ERSA may also be used for executing general-purpose applications that are less resilient to errors through the concept of configurable reliability. However, the associated costs can be higher.

Major contributions of this paper are:

1. Introduction of the concept of ERSA for probabilistic applications without requiring expensive error detection.
2. Detailed description of ERSA hardware and software architectures and probabilistic algorithm-aware optimizations.

3. Experimental results from ERSA hardware prototypes.
4. Analysis of computation accuracy and execution time trade-offs of ERSA over a wide range of error rates.

9.2 Error resilience of probabilistic applications

Our expectations on accuracy are different for different software applications. This is because applications have different degree of error tolerance in their algorithms. Probabilistic applications such as pattern recognition, data mining and case synthesis are good examples of applications with high error-tolerance [Dubey 05]. Probabilities are usually used for decision-making in these applications. As long as the low-order bit errors don't change the relative magnitude compared to the threshold used in decision-making, errors may not appear in the final result. Also, these applications reach final results after many iterations, providing opportunities for errors to be "averaged out" during iterations. ERSA bases on the idea that application algorithm level error resilience can be leveraged to lower the cost of hardware-level error resilience [Leem 10]. At the same time, algorithm-level error-tolerance cannot solely sustain high volume errors. Hardware architectural support is needed to protect system from types of errors that application algorithm is not error resilient to. Thus, ERSA is an excellent cross-layer reliability optimization between hardware system architecture and software application algorithm.

We use the following three probabilistic applications (from RMS benchmarks [Kestor 09] and related previous work [May 08]) to demonstrate the effectiveness of ERSA:

1. **K-means clustering:** A classification algorithm that partitions input data (points in n -dimensional space) into K clusters.
2. **Low-Density Parity-Check (LDPC) decoding:** A decoder module for the LDPC code, an error-correcting code that is widely used in communication applications. The decoding algorithm is based on loopy belief propagation.
3. **Bayesian network inference:** Bayesian network provides ways of "extracting and learning" new information from raw data. Given an image, our specific benchmark identifies cars in that image using its context (e.g., trees, roads, houses in the image).

9.3 ERSA Overview

ERSA is a multi-core architecture with asymmetric reliability (Fig. 9.1). ERSA consists of small number of highly reliable cores, referred to as *Super Reliable Cores* or *SRCs*, together with a large number of cores that are less reliable but account for most of the computation capacity, referred to as *Relaxed Reliability Cores* or *RRCs* (1 SRC and 8 RRCs in our hardware prototype). The motivation for asymmetric reliability comes from the computation model of probabilistic applications. An entire application can be divided into control-intensive resource management code (main thread in Fig. 9.2) that needs to be executed on error-free hardware while data-intensive computations are often more error-tolerant. By assigning the control-related code to SRC(s) and the computation intensive code to RRCs, we can minimize the number of processor cores that require high reliability and hence, avoid highly conservative overall system design. In our current ERSA implementation, interconnects between SRCs and RRCs are assumed to be reliable, e.g., using efficient error-correcting codes (ECC) [Raghunathan 03].

Two important points about ERSA are:

1. ERSA does not rely on traditional concurrent error detection in RRCs. Instead, it uses very few “light-weight” checks in RRCs, e.g., timeout and illegal memory access checks as discussed later.
2. ERSA is not restricted to probabilistic applications only. The ERSA platform may be

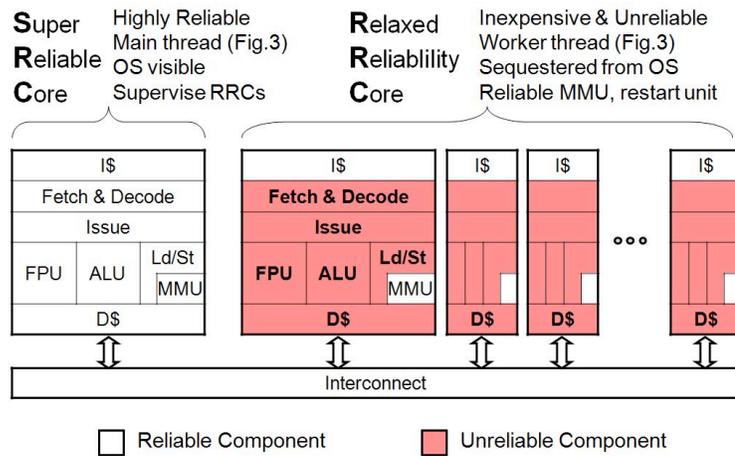


Figure 9.1. ERSA hardware architecture.

used to execute general-purpose applications as well using the concept of “configurable reliability” [Mitra 08]; i.e., according to application needs, reliability levels of various ERSA hardware components may be adjusted by modifying supply voltage, clock speed or by turning error protection mechanisms on/off.

9.3.1. Super Reliable Core (SRC)

An SRC is responsible for:

1. Executing the main thread that also performs RRC memory bounds check (Sec. 3.2) setup, RRC computation results reduction and convergence checking (Fig. 9.2).
2. Executing operations that are not resilient to errors (e.g., the operating system).
3. Distributing tasks to RRCs: ERSA uses task-queues to distribute tasks to RRCs during run-time. Our experience shows that run-time task distribution produces superior results than compile-time ones for ERSA. This is because not all RRC tasks complete execution in the presence of errors – RRCs can be non-responsive while executing tasks. Even when RRC tasks complete execution, they may not pass convergence checks (details in Sec. 3.3) performed by the SRC. All these tasks need to be dynamically reassigned to RRCs. Run-time scheduling with task-queues can adaptively reassign tasks to RRCs. It can also create “diversity” in the mapping of tasks onto RRCs (e.g., if a task repeatedly fails on a particular RRC, the run-time task scheduler can assign that task to another RRC).
4. Supervising RRCs: *Timeout bounds*, similar to watchdog timers [Mahmood 88], are

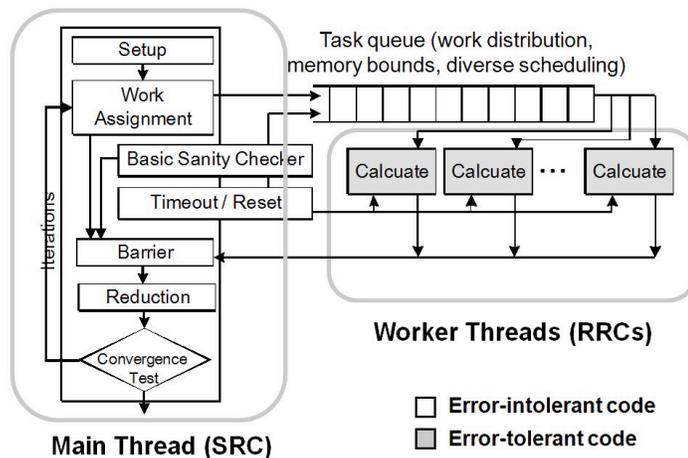


Figure 9.2. ERSA computation model.

used to detect when an RRC becomes non-responsive. The SRC checks for the liveness of an RRC via the corresponding completion bit in the task-queue. If it is not set within a timeout bound, the SRC terminates the execution of that RRC task and reboots it. In our ERSA experiments, we obtained timeout bound for each application through trial and error. In addition, SRC performs basic sanity checks on RRC computation results. Examples include probability value checks (must be between 0 and 1) and loop index checks (RRC loop indices are checked to determine if an RRC terminated early due to errors in loop indices).

9.3.2. Relaxed-Reliability Cores (RRCs)

RRCs constitute the large majority of on-chip processor cores in ERSA. They provide inexpensive and massive computing power. RRCs are sequestered from the operating system (OS) because the OS is highly vulnerable to errors [Kao 93]. RRC hardware is assumed to be unreliable except the Memory Management Unit (MMU) and the L1 instruction cache. The MMU is where memory access bound violations are detected. Memory bounds-checking is a popular way of detecting invalid memory accesses by producing additional checks for each read or write. ERSA uses memory bounds checks to prevent invalid memory write accesses from RRCs that may be caused by hardware errors. Memory bounds are obtained by collecting all the base addresses and size information of all static, heap and stack objects used in the RRC codes. Compiler-assisted approaches, e.g., [Jones 97], may be used for this purpose as well. In order to reduce the overhead of memory bounds checks, the same memory bounds are used for all memory instructions in an RRC. Incorrect memory accesses inside the memory bounds may create errors that are handled by ERSA similar to other error sources.

For the software code executed on an RRC, *function in-lining* is used to protect RRCs from stack pointer-related errors. Errors in stack pointer can overwrite the return address of a function on the stack so that, when the function returns, program control can jump to arbitrary locations. Most of the function calls in RRC codes are in-lined to protect the RRCs from the stack pointer errors at the expense of increased code size.

9.3.3. Algorithmic Convergence Test

The ERSA system described so far, i.e., incorporating asymmetric reliability and memory / timeout bounds checks and function in-lining is referred to as *Basic ERSA*. It is far more error-resilient compared to no-ERSA (details in Sec. 4). However, for error rates greater than 5,000 errors/sec/RRC or 5×10^{-5} error/cycle/RRC, *Basic ERSA* can result in large computation inaccuracies and significant execution time overheads. In order to overcome these challenges, algorithm-aware software techniques are needed. At high error rates, iterations in probabilistic applications often fail to converge, which leads to high execution time overheads. We designed a set of application algorithm-aware software techniques, which when used in conjunction with *Basic ERSA*, results in a system we refer to as *Enhanced ERSA*. Sections 3.3.1 and 3.3.2 present the relevant details.

9.3.3.1. Convergence Damping

Hardware errors in RRCs can result in fluctuating behaviors in computation results. Damping these fluctuations to make the computation results more stable is an effective way of mitigating the impact of hardware errors.

The damping scheme used in *Enhanced ERSA* is to “partially update” the output of an RRC; i.e., if the change (from the corresponding output in the previous iteration) exceeds a given amount, referred to as *saturation limit*, then simply update that output with the saturation limit. If the change is within the saturation limit, we simply accept the output of that RRC. In K-Means clustering, if a cluster diameter changes by more than +/-30% of the previous value (corresponding to a cluster assigned to an RRC), the coordinates of that cluster center are adjusted to make the diameter change by +/-30% of the previous value only. In LDPC decoding, the probability of a bit being ‘1’ is allowed to be updated by +/-0.25 at maximum until it reaches 0 or 1 in one iteration. Similarly, the probabilities of Bayesian-Network nodes are allowed to change by at most +/-0.3 per iteration. Choosing appropriate saturation limits is very important because over-damping can delay convergence (and increase execution time). For our ERSA experiments, we carefully chose these values through trial and error. One may choose not to use convergence

damping until the occurrence of a pre-determined number of iterations in order to accelerate convergence.

9.3.3.2. Convergence Filtering

Although convergence damping reduces fluctuations, it is better to discard the results from an iteration that shows excessively large fluctuations. *Convergence filtering* “selectively updates” the computation output by deciding whether to accept or to discard the results from all RRCs produced during the iteration. The difference between convergence filtering and damping is that filtering applies to the results produced by all RRCs during an iteration (i.e., all results accepted or discarded) whereas convergence damping is applied individually on the results obtained from each RRC. For example, LDPC decoding uses the total number of bits that failed parity check at the end of an iteration (referred to as *failed bits*) as the convergence criterion and declares convergence when the number of failed bits becomes zero. If this number increases by any amount, convergence filtering discards the results from that entire iteration. In K-Means clustering, the total sum of cluster diameters, which monotonically decreases to minimum value, is used as the metric for convergence filtering. If the sum increases by greater than 10% of the previous sum, the iteration is discarded. In Bayesian-network, the Euclidean distance between the probabilities of present and previous iterations is measured. If the distance is greater than 10% of the sum of probabilities in all nodes, new results are rejected. Similar to convergence damping, the parameters used for convergence filtering are obtained through trial and error. Convergence filtering can also be skipped until the occurrence of a pre-determined number of iterations.

9.4. ERSA Experiments

9.4.1. ERSA Experimental Results: Logic Errors in RRCs

We present experimental results from an ERSA prototype using PowerPC 405 100MHz cores in a Xilinx Virtex II Pro FPGA system (Fig. 9.3(a)). In our prototype system, we have one SRC and eight RRCs, which are mapped to two PowerPC hard cores in the FPGA. One PowerPC core is dedicated to the SRC and the other core is time-multiplexed to emulate RRCs. Time slices to each RRC process are assigned in equal

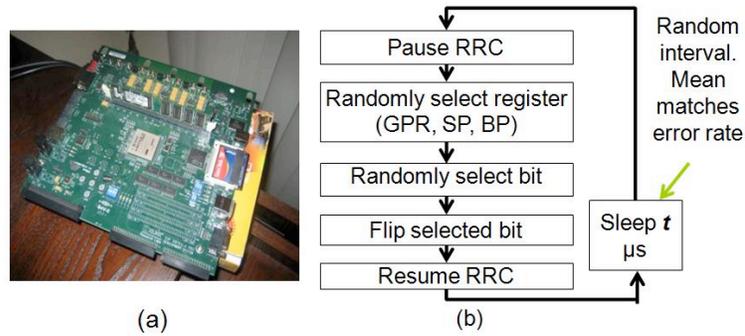


Figure 9.3. ERSA hardware prototype. (a) Xilinx Virtex II Pro board with Dual PowerPC 405. (b) Error injection model.

portions and in circular order. Separate stack memories are allocated for RRCs. For each emulated RRC, memory bounds checking is implemented using translation lookaside buffer (TLB) exception handler. At every TLB miss, memory access address is checked using the memory bounds (Sec. 3.2) assigned to the corresponding RRC. RRC timeout bound checking utilizes the hardware watchdog timer in PowerPC core, which times out and induces RRC core restart whenever its processor becomes non-responsive.

In this section, we demonstrate the effectiveness of ERSA through hardware error injections in RRCs. In the PowerPC 405 architecture, there are 32 general purpose registers (GPR) and no floating point registers. GPRs also serve as stack and base pointers. Thus, injecting errors into GPRs will emulate the effects of both control and data errors. RRC error injection flow is shown in Fig. 9.3(b). The error injection rate was varied from zero to 30,000 errors/sec/RRC, which corresponds to 3×10^{-4} error/cycle/RRC (processor clock frequency of 100MHz). In reality, hardware errors can be generated by any component such as adders, functional units, latches and flip-flops, etc. Injecting errors into those components directly will improve the accuracy of our experiments. This was not possible in PowerPC hard cores. Reconfigurable computing platforms may be used for more detailed experimentation.

The error injection routine is invoked by a programmable interrupt timer. One register from the 32 GPRs is randomly chosen and one bit out of 32 bit locations is randomly chosen and flipped. In PowerPC 405 architecture, there are other 32 special purpose registers that are used for controlling processor resources such as debugger and timers, etc. Since they are not visible to RRC code, errors were not injected into those registers.

Similar experimental setup was used in [Kao 93, Li 06, Wong 06] except that, unlike the ERSA hardware prototype, errors were injected using software simulators.

9.4.1.1. ERSA Computation Accuracy

We compare three implementations: *No ERSA*, *Basic ERSA* and *Enhanced ERSA* (Sec. 3.3). For the *No ERSA* case, probabilistic applications are executed under error injections without using the ERSA framework. *Basic ERSA* implements ERSA with memory/timeout bounds checks, task-queue, sanity checks and function in-lining (details in Sec.3). The *Enhanced ERSA* implementation includes *Basic ERSA* together with

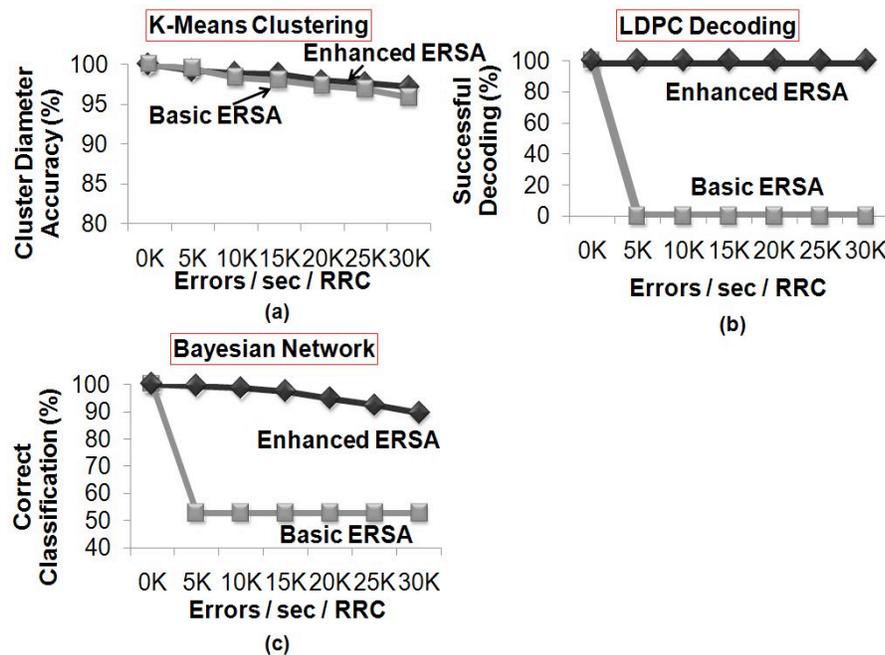


Figure 9.4. ERSA computation accuracy. Basic ERSA and Enhanced ERSA implementations are compared.

convergence damping / filtering. Figures 9.4(a)-(c) show the computation accuracy of ERSA. Individual metrics used to measure the computation accuracy are detailed in Table 1. Without ERSA (i.e., for the *No ERSA* case), the system crashes after 1 to 3 errors in average and no useful output is generated. *Basic ERSA* does not cause any system crash and achieves convergence. However, computation accuracy degrades significantly (by 50 to 100%) as error injection rate increases. K-Means clustering shows less degradation of computation accuracy, which can be explained by its high error resiliency due to averaging operations.

With *Enhanced ERSA*, computation accuracy achieved is better than 90% even for extremely high error rates of 30,000 errors/sec/RRC. In Fig. 9.5, we further analyze the cognitive resilience aspect of the Bayesian network application which exhibits 90% computation accuracy at 30,000 errors/sec/RRC (Fig. 9.4(c)). Bayesian network inference is used to identify cars in a satellite image. Objects inferred as cars are in green squares and others are in red. Depending on the final intended use of car identification, e.g., retrieving overall traffic conditions, 10% inaccuracy in Fig. 9.5 may not be very significant.

9.4.1.2. ERSA Execution time

Figure 9.6 shows execution times of ERSA applications. For *No ERSA*, execution times are virtually infinite because applications crash or do not converge even at very low error injection rates. With *Basic ERSA*, applications converge but execution time overheads can be very significant: 5 to 7 times compared to the no-error case. With *Enhanced ERSA*, the execution time overhead is less than 20 to 30% and increases gradually (shown in zoomed-in plots in Figs. 9.6(b,d,f) at extremely high error rates of 20,000 errors/sec/RRC. Note that, there is execution time overhead in ERSA even with zero errors. This overhead comes from convergence and sanity checks on RRC results. In the current ERSA implementation, RRCs may stall until the SRC completes the checking,

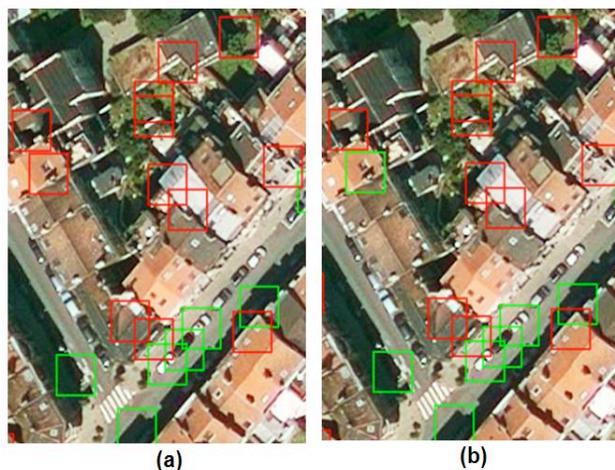


Figure 9.5. Output images of Bayesian Network Inference with (a) 100% (b) 90% accuracy. Objects inferred as cars are marked in green squares and others in red.

which can be improved in the future.

The execution time overhead of ERSA can be traded off with computation accuracy if necessary, i.e., the convergence criteria may be relaxed to make it easier to meet. Applications without errors can also expedite convergence using this technique [Chakradhar 09].

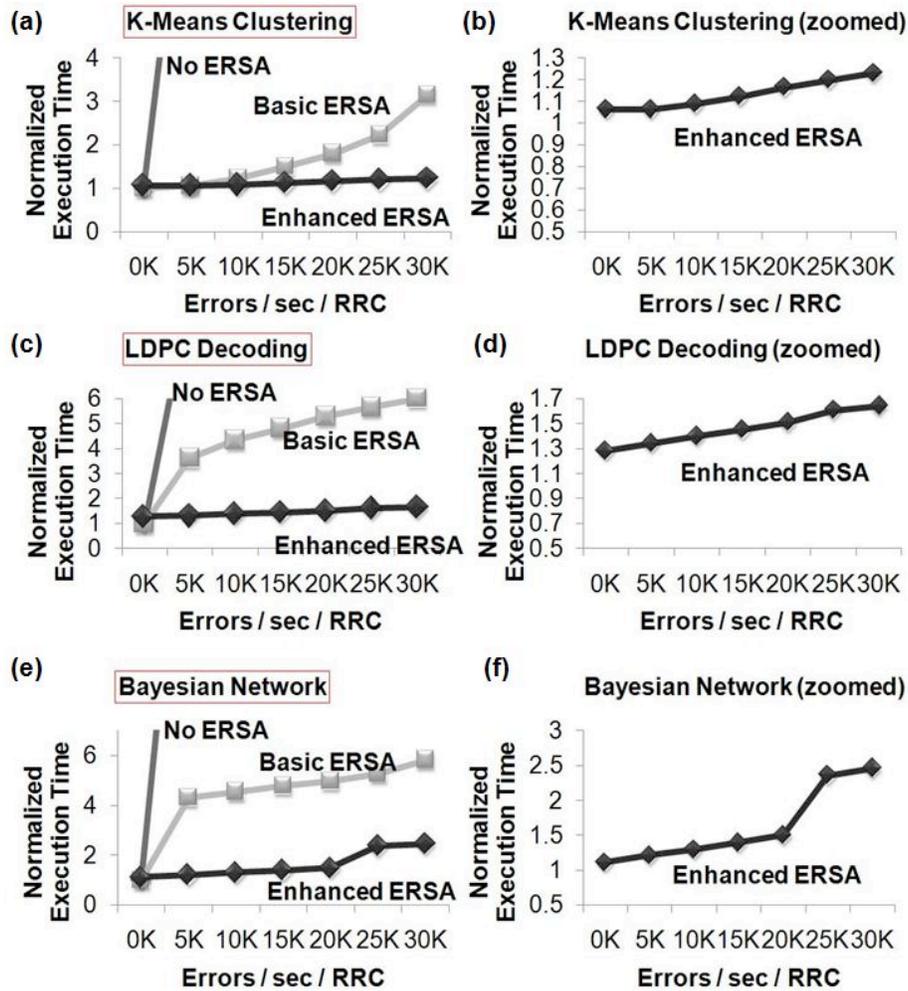


Figure 9.6. ERSA execution times (a, c, e) and the corresponding zoomed-in plots (b, d, f).

9.4.2. RRC L1 Data Cache Errors

ERSA techniques are also effective for non-transient static errors in SRAMs, e.g., in L1 data caches. SRAM errors are becoming significant in future technologies

[Mukhopadhyay 05, Agostinelli 05]. V_{ccmin} , the minimum voltage at which SRAMs can reliably operate, is a major challenge. SRAM V_{ccmin} errors have permanent locations because they originate from manufacturing-induced variations [Wilkerson 08]. There is another class of V_{ccmin} -related errors called erratic bit errors [Agostinelli 05] that are temporary in nature. Permanent locations of SRAM V_{ccmin} errors make probabilistic applications less error-resilient because the same errors continue to appear over iterations. ERSA overcomes this challenge by moving data around in the L1 cache (since error locations cannot be moved). From an application's perspective, moving data around has similar overall effects as changing error locations. This is accomplished by adding random offsets to memory addresses (Fig. 9.7(a)). In our ERSA implementation, each RRC has 16KB, two-way set associative data cache with 32 byte line size. A Linear Feedback Shift Register (LFSR) generates a random offset that is to be used over a given period of time. When data is stored in the cache, the offset is added to the address bits to move data to a new memory location chosen by the offset. This offset is subtracted from the address bits when the data is evicted, which makes using the address offset invisible outside the cache. New address offset is generated when there is an explicit erroneous event in the RRC such as RRC reboot or when RRC results are discarded during

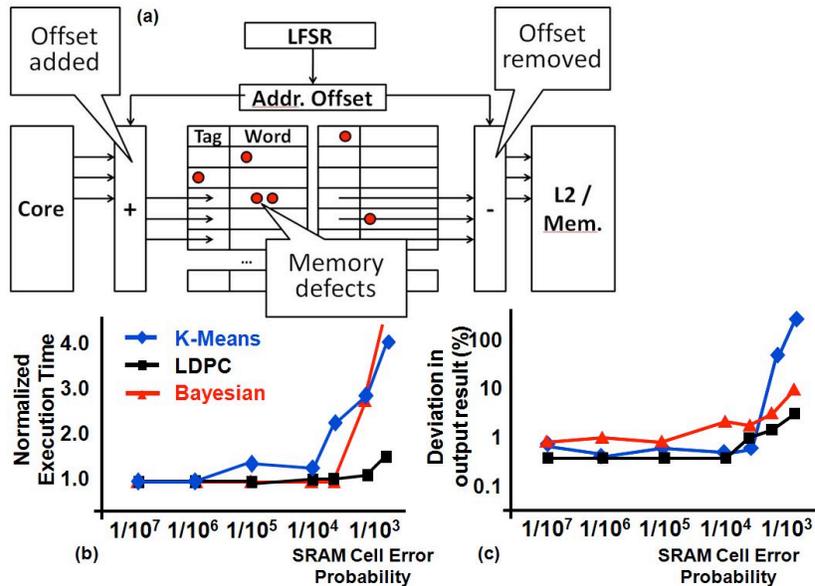


Figure 9.7. (a) ERSA L1 data cache organization, (b,c) data cache error experiment results.

convergence filtering.

As shown in Fig. 9.7, the execution time and computation inaccuracy gradually increase until they reach a “knee” point beyond which execution time and computation inaccuracy increase sharply (at SRAM bit error probabilities of 10^{-4} to 10^{-3} , which correspond to 0.32 to 3.2 % of erroneous cache words). Further research is required to push these knee points to higher error rates.

Reference

- [Agostinelli 05] M. Agostinelli, *et al.*, "Erratic fluctuations of SRAM cache Vmin at the 90nm process technology node", *Proc. Int'l Electron Devices Meeting*, pp.655-658, 2005.
- [Borkar 04] S. Borkar, *et al.*, "Design and reliability challenges in nanometer technologies", *Proc. Design Automation Conference*, pp.75, 2004.
- [Breuer 05] M. A. Breuer, "Multi-media Applications and Imprecise Computation", [*Proc. Euromicro Conference on Digital System Design*, pp.2-7, 2005.](#)
- [Chakrapani 06] L. Chakrapani, *et al.*, "Ultra Efficient Embedded SOC Architectures based on Probabilistic CMOS Technology", *Proc. Design Automation and Test in Europe*, pp.1110-1115, 2006.
- [Dubey 05] P. Dubey, "Recognition, Mining and Synthesis Moves Computers to the Era of Tera", *Technology at Intel Magazine*, 2005.
- [Eltawil 05] A. Eltawil, F. Kurdahi, "Improving Effective Yield Through Error Tolerant System Design", *Proc. Int'l Conference on Electronics, Circuits and Systems*, pp.125-129, 2005.
- [Gelsinger 06] P. Gelsinger, "Into the Core...", *Stanford Computer Systems Colloquium*, June 7, 2006.
- [Hayes 07] [J. Hayes](#), *et al.*, "An Analysis Framework for Transient-Error Tolerance", *Proc. VLSI Test Symposium*, pp. 249-255, 2007.
- [Huang 84] K. Huang and J. Abraham, "Algorithm-based fault tolerance for matrix operations", *IEEE Transactions on Computers*, Vol.C-33, Issue 6, 1984.
- [Jones 97] R.W.M. Jones, P. Kelly, "Backwards-compatible bounds checking for arrays and pointers in C programs", *Int'l Workshop on Automated Debugging*, pp. 13-16, 1997.
- [Kao 93] [W. Kao](#), *et al.*, "FINE: A Fault Injection and Monitoring Environment for Tracing the UNIX System Behavior under Faults". [*IEEE Trans. on Software Engineering*, Vol.19, Issue 11, pp. 1105-1118, 1993.](#)
- [Kestor 09] G. Kestor, *et al.*, "RMS-TM: A Transactional Memory Benchmark for Recognition, Mining and Synthesis Applications", *Proc. Workshop on Transactional Computing*, 2009.

- [Li 06] X. Li, D. Yeung, "Exploiting Soft Computing for Increased Fault Tolerance", *Proc. Workshop on Architectural Support for Gigascale Integration*, 2006.
- [Liu 91] J. Liu, *et al.*, "Algorithms for Scheduling Imprecise Computations", *Computer*, Vol.24, Issue 5, pp.58-69, 1991.
- [Mahmood 88] A. Mahmood, E.J. McCluskey, "Concurrent Error Detection Using Watchdog Processors-A Survey", *IEEE Trans. Computers*, Vol. C-37, No. 2, pp.160-174, 1988.
- [May 08] M. May, *et al.*, "A Case Study in Reliability-Aware Design: A Resilient LDPC Code Decoder", *Proc. Design Automation and Test in Europe*, 2008.
- [Mitra 05] S. Mitra, *et al.*, "Robust System Design with Built-In Soft-Error Resilience", *IEEE Computer*, Vol. 38, no. 2, pp. 43, 2005.
- [Mitra 08] S. Mitra, "Globally Optimized Robust Systems to Overcome Scaled CMOS Reliability Challenges", *Proc. Design Automation and Test in Europe*, pp. 941-946, 2008.
- [Mukhopadhyay 05] S. Mukhopadhyay, "Modeling of Failure Probability and Statistical Design of SRAM Array for Yield Enhancement in Nanoscaled CMOS", *IEEE Trans. CAD*, Vol. 24, No. 12, pp.1859-1880, 2005.
- [Pattabiraman 06] K. Pattabiraman, *et al.*, "Dynamic Derivation of Application-Specific Error Detectors and their Implementation in Hardware", *Proc. [European Dependable Computing Conference](#), pp.97-108, 2006.*
- [Rinard 03] M. Rinard, "Acceptability-oriented computing", *Conference on Object Oriented Programming Systems Languages and Applications*, pp.221-239, 2003.
- [Shanbhag 02] N. Shanbhag, "Reliable and energy-efficient digital signal processing", *Proc. Design Automation Conference*, [pp.830-835, 2002.](#)
- [Raghunathan 03] V. Raghunathan, *et al.*, "A survey of techniques for energy efficient on-chip communication", [Proc. Design Automation Conference, pp.900-905, 2003.](#)
- [Van Horn 05] J. Van Horn, "Towards Achieving Relentless Reliability Gains in a Server Marketplace of Teraflops, Laptops, Kilowatts, & "Cost, Cost, Cost"...", *Proc. Int'l Test Conference*, pp. 671-678, 2005.

- [Wilkerson 08] C. Wilkerson, *et al.*, "Trading off Cache Capacity for Reliability to Enable Low Voltage Operation," *Proc. Int'l Symp. on Computer Architecture*, pp. 203-214, 2008.
- [Varatkar 07] G. Varatkar, *et al.*, "Sensor Network-On-Chip," in *Proc. Int'l. Symp. on System-on-Chip*, 2007.
- [Wong 06] V. Wong, M. Horowitz, "Soft Error Resilience of Probabilistic Inference Applications", *Workshop on Silicon Errors in Logic-System Effects*, 2006.
- [Yu 00] S. Yu, *et al.*, "An ACS Robotic Control Algorithm with Fault Tolerant Capabilities", *Proc. Int'l Symp. Field-Programmable Custom Computing Machines*, pp.175, 2000.

Chapter 10.

Conclusions and future works

With all the fruits harvested from the low-power CMOS design tree, it has become clear that further improvements in energy efficiency will only be realized through a novel and revolutionary design solution [1]. That being said, we believe all three low-power research topics discussed in this thesis are disruptive and “out-of-the-box”. Novelties in these projects include following.

First, the Magnetic Coupled Spin-Torque Device (MCSTD) is a spintronic logic device that operates on magnetism rather than electrostatics. It is a non-transistor type logic device that has signal gain, fanout and signal level restoration capability, key ingredients for any logic technology. It can implement entire set of Boolean logic functions. Furthermore, it has an interconnection scheme that overcomes the short spin coherence length problem. At the time of writing this thesis, MCSTD is the only experimentally demonstrated spin-torque logic device that can implement the full set of Boolean logic and provides non-volatile logic functionality. The importance of MCSTD will grow for following reasons.

- i) The future of electronics will be dominated by the convergence of multi-functional devices: highly integrated, low-power and autonomous intelligent devices for healthcare, mobile and consumer applications. MCSTD serves as an excellent example for multi-functional logic+storage or logic+sensor+storage device.
- ii) Due to low-power techniques such as *power-gating* and *power-domains* [2], etc., logic portions of today’s embedded microprocessors do not consume



Figure 10.1 Sensory swarm. Trillions of simple devices spread in the environment adding sense to the internet [3]

power in stand-by mode. Therefore, memories have become the prime consumers of standby power (which is a major issue in mobile devices) [1]. Non-volatility combined with high speed and long endurance characteristics of MCSTD effectively eliminate the standby power consumed inside logic and memories altogether.

- iii) Lastly, and perhaps most importantly, with the proliferation of mobile and sensor devices, we will see more and more devices spread around in the environment, measuring things, sending data to the network, adding sensing to the internet and enriching our way to interact with information (Fig. 1). J. Rabaey of U.C. Berkeley calls this a “Sensory Swarm” [3]. It is expected that trillions of these devices will be deployed. These devices will soon become “truly immersive” and embedded into the environment and daily life. Potential applications are artificial skin [4], interactive surfaces [5], smart objects (smart tires, wearable computers), microscopic health monitoring, etc. However, current sophisticated power-saving techniques are not applicable in these devices. Also, many years of deployment time is required without changing batteries: they require another huge leap in size, cost and energy reduction.

MCSTD based logic makes a perfect solution for the “truly immersive” smart sensor applications: MCSTD does not require sophisticated power saving schemes due to its inherent non-volatility, whose zero standby power consumption will be an increasingly attractive low-power solution. The merged logic+storage+sensor features of MCSTD minimize communication delay and power among logic, I/Os and storage, which leads to ultra-low power architecture for smart sensors.

The second topic in this thesis, CNT/GNR heterojunction TFET, achieves a subthreshold slope less than 60mV/dec to make it possible for V_{dd} and V_{th} to further scaled down. This research effort can be considered as a disruptive low-power technique within the boundary of MOSFET-like devices. CNT/GNR heterojunction TFET completely eliminates the I-V ambipolarity problem for the first time. A type-II

heterojunction is achieved simply by partially unzipping the carbon nanotube. Strain-based carrier mobility engineering utilizes thermal relaxation in the Carbon atom bonds after unzipping the carbon nanotubes.

The third topic, Error Resilient System Architecture (ERSA) demonstrates a new paradigm for low-power system design. In the past decade, the state-of-the-art architecture-level low-power technique is to utilize concurrency to compensate for the slow-down in performance when clock frequency or supply voltage are reduced to save energy. ERSA is one of the first research efforts to examine the trades-off of reliability with energy efficiency. For example, scaling down supply voltage below V_{ccmin} will generate errors. ERSA gets additional robustness to mask out induced errors from 1) algorithmic resilience in software applications and 2) network of many-cores that can statistically sustain reasonable throughput. Actually, this idea is already adopted by nature, where the human brain consists of millions of computationally simple neurons, which have very low SNR. Low-power and robust computations are achieved by “*statistical computing*” performed by these neurons.

These three low-power ideas can be used in different design layers, i.e., device, circuit and architecture, to produce synergistic effects. With MCSTD gates used for general logic, we can finally achieve an “*instant-on computer*” and consume zero-static power. The output MTJ driver circuit is made of CNT/GNR heterojunction TFET that is more energy-efficient than MOSFETs. For system architecture, we choose ERSA architecture, which can trade-off reliability for low-power consumption. Despite the degradation in reliability, this system architecture is still reliable with the help of application-level error resiliency and network of many-cores robustness.

In this thesis, we also identified the shortcomings and weaknesses of these devices, which lead us to future works.

- i. MCSTD: First, the current version of MCSTD consumes more power than CMOS. This is because spin-torque transfer is not as energy efficient as needed: the switching current density is too large and the tunneling barrier is required for high TMR results in a large output resistance. We will need to investigate the option of using spin-valves for MCSTD gates and improve the

spin polarization and switching current of them. Heusler alloy based spin-valve seems to be a good candidate. Second, MCSTD requires CMOS circuits to supply switching current to the output MTJs. This CMOS intervention limits the maximum power saving that can be achieved. Incorporating magnetic domain-wall motion based interconnects and electric field induced magnetic reversal to MCSTD circuits will be useful in mitigating the power consumed in CMOS peripherals.

- ii. CNT/GNR heterojunction TFET: I_{on} did not increase simply by reducing the bandgap. It was limited by the density of states (DOS) of the source region.
- iii. ERSA: Overheads from reliability vs. energy trade-off increase faster than energy saving pays off.

None of these shortcomings are fundamental roadblocks. Considering the short history of low-power research in the individual areas discussed in this thesis, we believe that the results are very promising and project that there will be a rapid progress on these novel, disruptive technologies.

Reference

1. J. Rabaey, "Low Power Design Essentials," *Springer*, DOI 10.1007 (2009)
2. Royannez, et al., "90nm low leakage SoC design techniques for wireless applications," *Proc. of International Solid-State Circuits Conference* (2005)
3. J. Rabaey, "Statistical Computing : the alternative road to low energy," *Design Automation Conference* (2009)
4. T. Someya et al., "Conformable, flexible, large-area networks of pressure and thermal sensors with organic transistor active matrixes," *Proc. of National Academy of Sciences of the United States of America (PNAS)*, DOI:10.1073/pnas.0502392102 (2005)
5. F. Block et al., "VoodooSketch: extending interactive surfaces with adaptable interface palettes," *Proc. of the 2nd International Conference on Tangible and Embedded Interaction*, pp. 55-58 (2008)