Towards Wearout-Aware and Accelerated Self-Healing Digital Systems

A Dissertation

Presented to

the faculty of the School of Engineering and Applied Science

University of Virginia

in partial fulfillment of the requirements for the degree

Doctor of Philosophy

by

Xinfei Guo

May 2018

APPROVAL SHEET

This Dissertation is submitted in partial fulfillment of the requirements for the degree of Doctor of Philosophy

Author Signature:	Lufol
J –	

This Dissertation has been read and approved by the examining committee:

Advisor: Mircea Stan

Committee Member: Kevin Skadron

Committee Member: James Aylor

Committee Member: Samira Khan

Committee Member: Pradip Bose

Committee Member: _____

Accepted for the School of Engineering and Applied Science:

1PB

Craig H. Benson, School of Engineering and Applied Science

May 2018

In memory of my grandfather ... To my family and friends ...

> © Copyright May 2018 Xinfei Guo All Rights Reserved

Abstract

The down-scaling of CMOS technologies has continuously offered better performance, lower power and higher level of integration. However, in advanced nodes with smaller feature sizes, on-chip components such as transistors and interconnects are experiencing more aggressive degradations due to wearout (aging) effects, which are dominated by Bias temperature instability (BTI) and Electromigration (EM). Transistors become more susceptible to voltage stress due to the increased effective field with the scaling of the thin oxide. Similarly, the shrinking geometries of metal interconnects render higher current densities, and the tremendous number of transistors within a compact area has resulted in higher power densities as well. Together, these lead to increased on-chip temperatures which potentially accelerate the wearout effects. In the meanwhile, with the ubiquity of electronics (e.g. IoT) in our daily lives, there have been increasing demands for reliable system design. Many of such applications require very long lifetime, higher utilization rate and tighter hard error tolerances. They are possibly deployed in extreme environmental conditions, such as high temperatures, which, unfortunately, further accelerate wearout.

Conventional techniques of coping with wearout by "tolerating", "slowing down" or "compensating" still leave the wearout themselves unchecked since they keep accumulating fundamentally as the system operates. Moreover, the continuous increase of irreversible components of wearout over the entire lifetime will cause permanent errors and failures potentially. This thesis proposes and demonstrates a new category of techniques that "repair" wearout in a physical sense through accelerated and active recovery, by which wearout (both BTI and EM) can be reversed by actively applying several techniques, such as high temperature, negative voltages (for BTI) and reverse currents (for EM), thus leading to effective accelerated self-healing. By studying the frequency dependent behaviors of wearout and recovery experimentally, we demonstrate that the permanent portion of wearout can be almost fully eliminated and avoided by using in-time scheduled recovery. Our experiments on hardware demonstrate that the accelerated self-healing techniques for both BTI and EM wearout are fast, effective and feasible to implement.

To enable the on-chip implementations of the accelerated self-healing techniques and fully utilize the explored frequency dependent behaviors, we propose a cross-layer accelerated self-healing (CLASH) system which instruments the notion of recovery across multiple layers (from circuit to the system level) of the system stack. A full set of circuit IP blocks, including recovery boost components, novel BTI and EM sensors, multi-mode recovery assist scheme and novel power gating structures, is designed and implemented. As wearoutinduced failures become more visible at the system level, we also explore several potential architecture and system level solutions that are able to take advantages of intrinsic sleep behaviors for full recovery. Overall, these techniques work together to guarantee that the entire system performs for more of the time at higher levels of performance and power efficiency by fully exploiting the extra opportunities enabled by the accelerated self-healing.

Leading-edge nodes such as FinFETs endeavor to offer advantages of future scaled devices while offsetting the problems introduced by many generations of planar CMOS scaling. Adapting to the new challenges and fully benefiting from FinFETs require the growing knowledge and design experiences. To contribute to this knowledge base, in this thesis, we perform a comprehensive study based on circuit simulations across multiple technology nodes ranging from conventional bulk to advanced planar technology nodes such as Fully Depleted Silicon-on-Insulator (FDSOI), to FinFETs. As challenges such as wearout appear to be more pronounced in these advanced nodes, and this grow to be critical especially in IoT applications in which industry is in the process of updating technologies. Each of those end markets has unique needs and characteristics, which affects how chips are used and under what conditions. We investigate how wearout can impact different categories of future loT applications with foundry-provided wearout models. We conclude that wearout needs to be considered in the full design cycle and the IoT lifetime estimation requires to incorporate wearout as an important factor. IoT-specific design solutions for mitigating wearout are also presented in this thesis.

Acknowledgements

It has been my most memorable experiences in the last five and half years, this journey would not have been possible without the help from many people.

First and foremost, I am deeply grateful to my advisor Prof. Mircea Stan, for his guidance, discussions and endless support both academically and financially. He is one of the smartest researchers I have ever met, he taught me to think out of boxes while still focusing on details. I'd like to thank him for offering me complete freedom to explore my own research ideas. I also owe a substantial amount of gratitude to his generous support for attending many international conferences, presenting my work to the outside world and applying for awards. It was from these experiences I learned a ton from peers and this helped me shape the research directions.

I wish to thank my Ph.D. committee members: Prof. Kevin Skadron, Prof. James Ayler, Prof. Samira Khan, Prof. John Lach and Dr. Pradip Bose for their valuable suggestions and comments. During the last year and half, it was my fortune to work directly with Dr. Bose and his team at IBM research, it was through this close collaboration and those weekly discussions I learned how it looks like in an industry research team and how to be an organized professional. I also want to thank Prof. Skadron and his former research associate Dr. Mohamed El-Hadedy (now an Assistant Professor at CalPoly) for offering me opportunities to work together on many interesting projects which all resulted good publications. I appreciate their time for writing me reference letters as well. A special thanks to Prof. Wayne Burleson for useful feedback and discussions on guiding early directions for this thesis.

I had the most fortunate experiences to work closely with colleagues from my lab and other research groups. They are Alec Roelke, Vaibhav Verma, Patricia Gonzalez-Guerrero, Sergiu Masanu, Dr. Kaushik Mazumdar, Dr. Divya Akella Kamakshi, Dr. Linqiang Luo and Dr. Harsh Patel. I learned a lot from them, and I thank them for their contributions to this thesis. Other former and current members and visitors of the HPLP family have also played vital roles in my life at UVa, including Mandi Das, Yunfei Gu, Junhan Han, Oluwole Jaiyeoba, Dr. Mehdi Kabir, Ben Melton, Mateja Putic, Tommy Tracy II, Elena Weinberg, Andy Whetzel, Theodoric Xie, Dr. Yingbo Zhao and Dr. Matthew Ziegler. I am grateful to

our system administrator Gary Li and department staff members Natalie Edward and Terry Tigner for their patience and help through my Ph.D. journey.

My graduate life wouldn't be so eventful and enjoyable without friends in Charlottesville, including Dr. He Qi, Hang Zhang, Dr. Runjie Zhang, Dr. Wei Zhang, Luonan Wang, Ningxi Liu, Kapila Wijayaratne, Abbis Haider, Dr. Yanqing Zhang, Dr. Ke Wang, Minyao Zhang and Dr. Liang Wang. It was my pleasure to know them and share my journey with them. My long-time friends Feng Gao, Yang Yu, Guoqing Fu and Yong Wang have always been there when I needed, thanks for their encouragement and support.

Last but certainly not least, I want to thank my family for their unconditional love and support of my dreams. As a small-town boy from northern China, I could have never achieved what I have now without the support from my parents, who always offered me the best they can. My Mom was also my first English teacher, and she always ensured I had the best resources to learn English. My father has been my math and physics tutor since middle school. Both of them are my lifelong teachers for all aspects of life, and they set me the best examples of being a good person. As the only kid in family, I owe them a lot for being away from home for so many years. My grandfather has been a mentor, a role model and a great friend for me, I still remember the time we spent together talking about sports, calligraphy and everything in life. His rich experiences, broad knowledge and optimism have impacted me endlessly. I know he is looking down on me from heaven and hope I can become someone he is always proud of. It was the best thing in my life to meet my girlfriend Jialin Luan, she has always been understandable, patient and supportive. We have so many similarities, and we share passions and dreams about our life. It was her love through the last one and half years that kept me motivated and made me stronger. I can't wait to share the next journey with her to our bright future.

There is an old Chinese saying "A drop of water shall be returned with a burst of spring". I'd like to dedicate this thesis to everyone who helped me at different stages of my life.

This work was supported by an IEEE CASS Pre-doctoral Fellowship, by DARPA under the PERFECT program and the UPSIDE program, by NSF grant CCF-1255907, SRC task 2410.001, and by C-FAR, one of six SRC STARnet Centers, sponsored by MARCO and DARPA.

Table of Contents

Li	st of H	figures		XV
Li	st of 7	Tables		XXV
1	Intro	roduction		
	1.1	Motiva	ation and Background	. 1
	1.2	Wearo	ut Mitigation Techniques	. 5
	1.3	Thesis	Contributions	. 8
	1.4	Thesis	Organization	. 11
2	Acce	elerated	Self-Healing Techniques for BTI Wearout	15
	2.1	Overvi	iew	. 15
	2.2	BTI W	Vearout and Recovery Basics	. 17
	2.3	Prior V	Work on BTI Recovery	. 20
	2.4	Accele	erated Self-Healing	. 22
		2.4.1	Active and Accelerate BTI Recovery	. 22
		2.4.2	Gate-level Analytical Model for Accelerated Self-Healing	. 22
	2.5	Experi	mental Setup	. 24
		2.5.1	Test Platform	. 24
		2.5.2	Test Configuration	. 26
		2.5.3	Test Flow and Test Conditions	. 27
		2.5.4	Modeling BTI Stress and Recovery for FPGA Test Structures	31
	2.6	Test R	esults for Accelerated BTI Wearout	. 34
		2.6.1	AC Stress vs. DC Stress	. 34
		2.6.2	Effect of Temperature on BTI Wearout	. 35
	2.7	Test R	esults for Accelerated Self-Healing Techniques	36
		2.7.1	Negative Voltage	. 37
		2.7.2	High Temperature	. 37

		2.7.3	Model Validation	40
	2.8	Revers	ible vs. Irreversible BTI Wearout	41
		2.8.1	Fast Traps vs. Slow Traps — A Physics Perspective	42
		2.8.2	Irreversible Wearout during Accelerated Self-healing	44
		2.8.3	Sequentiality of Reversible and Irreversible Wearout	45
	2.9	Freque	ncy Dependency of BTI Wearout and Recovery	47
		2.9.1	Sleep with Accelerated Rejuvenation when Getting Tired	47
		2.9.2	Measurement Results	50
		2.9.3	Reduction of Necessary Design Margin	53
		2.9.4	Reduction of Tracking Power	54
		2.9.5	Average Performance Improvement	55
		2.9.6	Frequency Dependency Behaviors of BTI Wearout	57
	2.10	Conclu	sions	58
	2.11	Acknow	wledgments	59
3	Acce	lerating	and Activating Recovery for EM Wearout	61
C	3.1	Overvi	ew	61
	3.2	EM We	earout and Recovery Mechanisms	63
	3.3	Prior W	Vork on EM Recovery	65
	3.4	"Rever	sing" the Direction of EM Wearout	66
	3.5	Test Se	etup	67
		3.5.1		67
		3.5.2	Measurement Setup	69
		3.5.3	Test Cases	70
	3.6	Experin	mental Results for EM Active and Accelerated Recovery	71
	3.7	EM Sig	gnoff Considering Accelerated and Active Recovery	75
		3.7.1	Relax the EM Design Rules	75
		3.7.2	Performance Improvement	79
		3.7.3	Extend the Wire Lifetime	80
	3.8	Summa	ary: EM vs. BTI	81
	3.9	Conclu	sions	82
	3.10	Acknow	wledgements	83
4	Circ	uit Tech	iniques for Accelerated and Active Recovery	85
	4.1	Overvi	ew	85
	4.2	Circuit	Solutions for Activating and Accelerating BTI Recovery	87
		4.2.1	On-Chip Negative Voltage Generation	87

		4.2.2	Negative Bias Voltage in a Logic Path	90
		4.2.3	Wearout-aware Power Gating	95
		4.2.4	On-Chip Heat Generation	100
	4.3	Circuit	Solutions for Activating and Accelerating EM Recovery	105
	4.4	BTI Se	nsing	111
		4.4.1	Previous BTI Sensing Techniques	112
		4.4.2	Ring Oscillator-based Test Structures for Separating NBTI and PBTI	113
		4.4.3	Metastable-Element-based Embeddable BTI Sensors	115
	4.5	EM Set	nsing	130
	4.6	Conclu	sions	135
	4.7	Acknow	wledgements	136
5	Acce	elerated	Self-Healing as a Key Design Knob for Cross-Layer Resilience	137
	5.1	Overvi	ew	137
	5.2	Accele	rated and Active Recovery Space Exploration	139
		5.2.1	High Temperature or Negative Voltage? or Both?	139
		5.2.2	Right Balance of Wearout and Recovery for BTI	141
		5.2.3	Right Balance of Wearout and Recovery for EM	142
	5.3	Archite	ecture-level Accelerated Self-Healing	143
		5.3.1	Architectural-level Model for Wearout and Lifetime Analysis	143
		5.3.2	Unit-level Accelerated Self-Healing	144
		5.3.3	Dark Silicon and Core Redundancy	146
	5.4	Schedu	ling at the System Level	148
		5.4.1	Reactive Recovery vs. AC Stress	148
		5.4.2	Application-dependent Proactive Recovery with Scheduling	150
	5.5	Recove	ery-Driven Design Methodology for Resilient System Design	153
	5.6	Putting	It All Together – CLASH: Cross-layer Accelerated Self-Healing	
		System		155
	5.7	Tradeo	ff Analysis	157
	5.8	Conclu	sions	159
	5.9	Acknow	wledgements	160
6	Expl	oring C	Circuit Aging in FinFET-enabled Internet of Things (IoT) Applica	-
	tions	5		161
	6.1	Overvi	ew	161
	6.2	Back to	the Future: Digital Circuit Design in the FinFET Era	163
		6.2.1	Motivation for Studying FinFET Technology	163

		6.2.2	FinFET Scaling and Sizing	165
		6.2.3	A Comprehensive Study of Bulk vs. FDSOI vs. FinFET Devices	168
		6.2.4	FinFET Fabrication	177
		6.2.5	FinFET Circuits	180
		6.2.6	FinFET Technology for Energy-constrained IoT Applications	191
		6.2.7	Summary - Digital Circuit Design with FinFETs	193
	6.3	When	"things" get older – Exploring Transistor Aging in IoT Applications .	195
		6.3.1	Motivation	195
		6.3.2	Previous Work on Aging in IoT Domain	196
		6.3.3	IoT Application Domains	198
		6.3.4	Simulation Results	201
		6.3.5	IoT Lifetime: Battery vs. Chip Lifetime	207
		6.3.6	Potential Solutions for IoT Circuit Aging	209
	6.4	Conclu	isions	213
	6.5	Ackno	wledgements	214
7	Con	clusion	s and Future Directions	215
	7.1	Summ	ary of Contributions	215
	7.2	Future	Directions	217
		7.2.1	Accelerated Self-Healing in Emerging Technologies	218
		7.2.2	Exploring Other Sources for Accelerating Recovery	218
		7.2.3	Integrating Wearout and Recovery in EDA Design Flow	219
		7.2.4	Dynamic Wearout Management by Self-learning	219
		7.2.5	Teaching Wearout and Recovery as Part of the VLSI Classes	220
Aj	opend	ix A L	ist of Publications	221
	A.1	Peer-R	eviewed Journals	221
	A.2	Peer-R	eviewed Conferences and Workshops	221
	A.3	Poster	8	223
	A.4	Presen	tations/Talks	224
Aı	opend	ix B F	low for Placing Metastable-element-based BTI Sensors	225
	B.1	New T	op-down Design Flow with BTI Sensor Insertion	225
	B.2	Demo	nstration in a Counter Design	226
Aj	opend	ix C A	Brief Overview of My Side Projects	229
_	C.1	Post-si	licon Hold Time Closure – Tunable Buffer Insertion	229

Bibliography		
Glossary		235
C.5	On-chip Power Regulation with Voltage Stacking	232
C.4	Dual-Data Rate Transpose Memory	231
C.3	Programmable Processing Element for IoT Crypto Systems	230
C.2	A 14nm Low-Vdd Heterogeneous RISC-V-based SoC	230

List of Figures

1.1	Wearout vs. technology scaling projected by Intel [151]. Y-axis refers to	
	wearout related metric, and it is normalized to the 32nm node	2
1.2	Illustration of BTI and EM. The Scanning Electron Microscope (SEM) figure	
	(from [210]) shows a cross-section of a fully processed microprocessor. BTI	
	occurs to transistors, and EM happens in metal wires.	4
1.3	Taxonomy of BTI and EM Mitigation Techniques. Note that "Recovery	
	Boost" here refers to the existing recovery acceleration solutions which	
	mainly focused on BTI wearout for SRAMs	5
1.4	Thesis Organization.	12
2.1	BTI Stress and Passive Recovery: NBTI happens in PMOS transistors; PBTI	
	happens in NMOS transistors. For both mechanisms, BTI starts recovering	
	when transistors are turned OFF, but this passive recovery period is very slow	
	and unpredictable.	16
2.2	Two BTI Mechanisms.	17
2.3	BTI behaviors modeled by Trapping/Detrapping (TD) theory. Passive Re-	
	covery still leaves a net ΔV_{th} that is almost permanent and hardly recovered	
	within a reasonable time.	19
2.4	Proposed accelerated self-healing solutions for NBTI. Similar solutions can	
	be applied to PBTI as well.	21
2.5	Illustration of a potential use case flow for Gate-level Accelerated Self-	
	healing model. This flow can be used for evaluating how recovery conditions	
	can affect the circuit metrics and system lifetime further	24
2.6	FPGA test platform (Lattice Semiconductor iCE40 HX-Series) architecture	
	[89]	25
2.7	BTI Test Configuration on FPGAs.	26

2.8	BTI Test Flow: FPGA chip communicates with computer through an Atmel	
	evaluation board, on which a micro-controller can be programmed with C.	
	The inherent timer in MCU can be programmed to controller how long the	
	chip is stressed or recovered.	28
2.9	BTI Test Setup: FPGA chip is placed in a socket board, which connects with	
	the interface board through a flat cable to separate the micro-controller from	
	the high temperature environment, only the FPGA chip is in thermal chamber.	30
2.10	Pass-transistor based LUT structure.	33
2.11	AC and DC Stress Measurement Results: AC stress case with a 50% duty	
	cycle shows slower BTI wearout.	35
2.12	Accelerated BTI Wearout under 110°C and 100°C for 1 day	36
2.13	Negative voltage-enabled active recovery after being stressed for 24 hours	
	(Net delay increase is $\sim 3.24ns$). X-axis is the recovery time. (a) at 20°C,	
	(b) at 110°C	38
2.14	High Temperature-accelerated recovery after being stressed for 24 hours	
	(Net delay increase is $\sim 3.24ns$). X-axis is the recovery time. (a) under 0V	
	(<i>passive</i> recovery), (b) under $-0.3V$ (<i>active</i> recovery)	39
2.15	Delay change (ΔT_d) over time during recovery: Model predictions vs. Mea-	
	surement. Solid lines are model predictions, markers are measured data	41
2.16	One test case showing that accelerated self-healing techniques can recover	
	about 72.4% of BTI wearout within only 1/4 of the stress time. Passive	
	recovery data are not shown in the figure, but it is about less than 2% of	
	recovery after 12 hours.	41
2.17	Energy distribution of electrons at room temperature	43
2.18	Illustration of Fast traps vs. Slow traps.	44
2.19	Energy distribution of electrons at room temperature.	45
2.20	Irreversible part under two recovery conditions – Passive recover vs. Ac-	
	celerated & active recovery (accelerated self-healing). BTI doesn't recover	
	completely even under the accelerated self-healing conditions	46
2.21	Accumulation of the irreversible BTI wearout cycle by cycle. IRn refers to	
	the accumulated irreversible component after nth cycle	46
2.22	BTI-induced frequency degradation under two accelerated stress conditions.	
	In both cases, reversible wearout kicks in firstly, then it levels off and irre-	
	versible wearout takes over	48
2.23	Sequentiality of reversible and irreversible wearout.	48
2.24	Measurement results for different "circadian rhythms"	51

2.25	Irreversible component accumulated during the first 6 cycles for four different	
	scheduled stress and accelerated and active recovery shifts (circadian rhythms).	52
2.26	An identical regular-operation use case (31 hrs vs. 1 hr) to the 1 hr vs. 1 hr	
	accelerated stress case for FULL recovery	53
2.27	Necessary design margin estimation under different stress conditions	54
2.28	Average performance improvement (IMP) for 1 day and 2 days from mea-	
	surement.	56
2.29	Average performance under different stress conditions	56
2.30	Frequency dependency of irreversible component of BTI wearout	57
2.31	Chapter 2 highlights.	59
3.1	Illustration of EM Recovery: Stress relaxation occurs when current is	
	switched off, and this is similar to the passive recovery condition for BTI.	62
3.2	Electromigration Mechanism: EM is the result of the dominant force F_{wind} ,	
	that is, the momentum transfer from the electrons which move in the applied	
	electric field.	63
3.3	Electromigration Stress and Passive Recovery: EM mainly affects the power	
	delivery network (PDN). EM healing occurs when current is removed, but	
	the recovery is partial and slow	64
3.4	Electromigration "reversing" techniques: No. 1 refers to the passive recovery,	
	No. 2, 3 and 4 are proposed active and accelerated recovery solutions	67
3.5	Die photo with the test structure for EM recovery: On-chip "long" and	
	"narrow" metal lines and their dimensions; "rt" means room temperature	
	$(\sim 27^{\circ}C)$	68
3.6	Electromigration stress and recovery measurement setup	69
3.7	Measurement results for EM degradation and recovery under passive recovery	
	(Figure 3.4 No. 1 or Test case EMPR1 in Table 3.1) and proposed recovery	
	conditions (Figure 3.4 No. 4, at 230°C and $\pm 7.96MA/cm^2$) during the void	
	growth phase: there is still a permanent component even under accelerated	
	and active recovery	72
3.8	Measurement results for EM accelerated and active recovery during the early	
	period of the void growth phase (at 230°C and $\pm 7.96MA/cm^2$): full recovery.	73
3.9	Measurement results for scheduled periodic recovery intervals during void	
	nucleation phase: It takes much longer for voids to nucleate, and the overall	
	TTF is extended significantly	74
3.10	Illustration of EM Current Rule Relaxation due to Recovery: Current density	
	requirement can be potentially relaxed, x is a number larger than 1	76

3.11	Illustration of cases with and without recovery during normal operation: If the current limit is doubled, the metal wire ages faster (almost twice compared to before), but the in-time accelerated and active recovery can always bring it back to fresh state. Overall, the EM current limit can be	
	relaxed while assuring the reliable operations during the lifetime span	79
3.12	Illustration of IR drop on power mesh	80
3.13	Chapter 3 highlights.	83
4.1 4.2	Overview of Chapter 4	86
4.3	A switch-cap based negative voltage generator (designed in 28nm FD-SOI technology) for delivering the negative voltage for BTI active recovery. <i>Vout</i>	88
4.4	is the negative voltage output	88
4.5	$315um^2$) is only ~ 5% of the total PMU area	89
	refers to the node voltage at node x , $n0$ is equal to $0V$	90
4.6	An equivalent circuit for ease of analyzing the behaviors under active BTI	~ .
4.7	recovery enabled by negative voltage supply $-V_{rec}$	91
4.8	slowly (an example of recovery voltage at $-0.3V$ shows this) PMOS Transistor V_{gs} for each logic stage under different negative voltage supply levels. The larger V_{gs} is, the more negative bias PMOS transistors experience, the higher BTI recovery rate will be. As the logic goes deeper,	91
4.9	V_{gs} decreases slowly. Depending on the logic depth, negative voltage level needs to be carefully picked	93
	nominal voltage, $V_{in} = 0$ or V_{dd}).	94

4.10	Power gating structure: it can be used effectively to reduce leakage and	
	help BTI recover passively, retention registers are used for saving the states.	
	The voltage drop due to the resistance of header transistor ΔR can lead to	
	performance loss, it is even worse that this transistor is ON most of the time	
	and experience BTI wearout as well.	95
4.11	Sleep transistor threshold voltage increase ΔV_{th} vs. Load performance loss	
	(8 49-Ring Oscillators running in parallel, the header transistor is sized as	
	$1\mu m$ wide)	96
4.12	An active recovery-enabled power gating structure: Blue outline parts refer	
	to the additional logic on top of the existing infrastructure. Sleep signal is	
	the trigger for starting the active recovery	97
4.13	Functional simulation of the proposed wearout-aware power gating structure	
	in 14nm bulk FinFET technology, signal names correspond to the ones in	
	Figure 4.12: When <i>Sleep</i> signal is high, the BTI active recovery mode starts,	
	the negative voltage is delivered as the virtual supply for the logic, the header	
	sleep transistor is also in reverse biased because a higher-than- V_{dd} voltage	
	is applied to the gate. The switching time between modes are similar to the	
	ones in existing power gating infrastructure	98
4.14	The physical layout of the power gating that implements the BTI active	
	recovery logic in 28nm FD-SOI technology. The added overhead is only	
	a NMOS header and some control logics. The load is 8 ring oscillators in	
	parallel, and the frequency divider is used for off-chip frequency readout.	99
4.15	Reconfigurable On-chip Heating Elements Schematic: The output frequency	
	of the ring oscillators can be reconfigured by selecting the length of the	
	inverter chain $(L, L/2, \text{ etc.})$. Different frequencies of the heating elements	
	correspond to a wide range of temperatures	101
4.16	Maximum temperatures that correspond to different oscillation frequencies	
	with the heating elements on FPGAs, temperatures are sampled with external	
	thermal sensors, a precalibration has to be carried out on the target FPGA to	
	determine this relation, data is from [228]	102
4.17	Potential strategies of on-chip reconfigurable heaters: (a) Evenly distribution;	
	(b) Ring placement; (c) Critical path placement. Red squares refer to the	
	heating element. An external controller is needed to configure the heater and	
	select the blocks that need to be recovered	103

4.18	Assist circuitry for activating BTI and EM recovery: (a) The main circuitry,	
	arrows represent the current direction under different modes, V_{DD} and V_{SS}	
	pins can be connected to the on-chip voltage regulator directly, or to the	
	global power delivery network; (b) Truth table for three operating modes;	
	(c) An example of activating NBTI recovery under BTI Active Recovery	
	mode, for PBTI recovery, the input needs to be "0", ΔV represents voltage	
	droop/increase or noise.	106
4.19	Three modes of the assist circuitry for activating BTI and EM recovery:	
	Normal Operation, EM Active Recovery and BTI Active Recovery. VDD and	
	VSS grid in the figure refer to power delivery network (PDN)	107
4.20	Functionality Simulation for EM/BTI Active Recovery Assist Circuitry in	
	28nm FD-SOI technology	108
4.21	Load Size vs. Performance and Switching Time: Increasing the number of	
	loads will reduce the performance and increase the switching time between	
	modes, to compensate the degradation, header/footer transistors need to be	
	upsized, which will further increase the area. This tradeoff needs to be	
	carefully considered during the design process.	109
4.22	Vertical cross section of the physical implementation for the EM/BTI Active	
	Recovery Assist Circuitry for VDD Grid (VSS Grid can be implemented in	
	a similar way): EM hazards happen at high current density regions, which	
	could be caused by faster switching activities on the load logic; At the logic	
	level, BTI hazards happen due to the continuous stress.	110
4.23	Sensing element for separating NBTI and PBTI. (a) Functional schematic	
	of 1 sensing stage. Red refers to the actual logic that experiences N/PBTI.	
	Odd number of stages can form a ring oscillator and serve as test structures	
	or sensors. (b) Truth table for the sensing stage. When under NBTI or PBTI	
	mode, PMOS or NMOS transistors are under constant stress; when under	
	Test mode, the structure works as a regular ring oscillator stage, and the	
	oscillation frequency can be read out from the RO output	114
4.24	Layout of an example of using sensing element in a 37-stage ring oscillator	
	in 28nm FD-SOI technology. Overall area is small, and the sensor can be	
	configured to sense either NBTI or PBTI. It can also be used as the test	
	structure for studying N/PBTI degradations.	114
4.25	Illustration of the BTI-induced Critical Path Reranking	116
4.26	Multiple-Critical-Path Embeddable BTI Sensor high-level scheme. (a) Proac-	
	tive Recovery Sensing; (b) Multiple-Paths Sensing.	118

4.27	Transistor-level schematic of the BTI sensor	119
4.28	Functionality simulation of the BTI sensor under fresh ($time = 0$) condition.	120
4.29	Simulation setup for introducing BTI wearout in a circuit netlist	121
4.30	Threshold degradation and sensor trigger point	122
4.31	Sensor triggering order for the proactive recovery case	123
4.32	Input patterns for multiple-path simulation case	124
4.33	Sensor trigger order in Multiple Critical Paths case.	125
4.34	Monte Carlo Simulations at the "fresh" status (500 points). More than 85%	
	of the seeds give the correct outputs. The overlapping color in the figure	
	refer to waveforms for each seed. The output of the sensor at TT corner is	
	used as the reference	126
4.35	Layout of a 2-path version of the proposed BTI sensor in 28nm FD-SOI	
	technology. The total area is only $54.1um^2$.	126
4.36	A 2-path version of the sensor implemented with the PnR tool directly	127
4.37	Embedded sensor in a scan chain cell in support of both close-loop and	
	open-loop sensor readouts	128
4.38	New scancell design methodology with the BTI sensor embedded in a Syn-	
	opsys design environmentCEL and .FRAM are layout formats required by	
	the Synopsys tools. Similar methodology can also be adapted to Cadence	
	design environment	129
4.39	Illustration of Metal-line-based EM sensors. Multiple dimensions can be	
	used to sense at different levels.	131
4.40	The EM-induced resistance change detection circuit (design is modified	
	based on the circuit proposed in [77]).	132
4.41	Illustrations of EM sensor usage in the accelerated and active recovery case.	133
4.42	Chapter 4 Highlights	135
5.1	Illustration of Concept of Cross-Layer Techniques.	138
5.2	Accelerated self-healing space exploration for a case when a 6-hour recovery	
	follows a 24-hour accelerated stress. Accelerated Recovery - high tempera-	
	ture, Active Recovery - Negative voltage.	140
5.3	Illustration of the "training" process for finding the optimal stress vs. recov-	
	ery balance to fully recover from the BTI wearout effects	141
5.4	Illustration of finding the optimal stress vs. recovery balance for EM during	
	run time	142

5.5	Heat (a) and Wearout (b) maps representing the average temperature and	
	relative wearout rate, respectively, when running <i>cholesky</i> [23]. In (a), red	
	indicates hotter temperatures and blue indicates cooler ones. In (b), red	
	indicates faster wearout while blue indicates slower wearout	145
5.6	A potential self-healing solution in a multi-core system. Dark silicon and	
	core redundancy can be utilized to improve the lifetime of the whole system.	
	t1 and $t2$ are two time points during the lifetime	147
5.7	Illustration of Reactive Recovery vs. Proactive Recovery vs. AC Stress	149
5.8	Recovery time under different accelerated & active recovery conditions after	
	12-hour constant stress under regular operation condition (room temperature,	
	nominal V_{dd}).	151
5.9	An illustration of recovery-driven design methodology in an IC design cycle.	153
5.10	An illustration of Cross-Layer Accelerated Self-Healing (CLASH) System.	155
5.11	Illustration of periodic proactively scheduled EM/BTI recovery.	156
5.12	Cross-layer implementation of accelerated self-healing for leveraging the	
	tradeoffs. At the circuit level, recovery circuit and wearout sensors are	
	distributed on the wearout-critical units, and they are triggered when needed.	
	Architecture level accelerated self-healing solutions can utilize some intrinsic	
	sleep behaviors and heat to recover inactive parts. It can also compensate	
	some of the power overhead introduced by the circuit level recovery solutions.	
	System level scheduling is able to divide the recovery tasks and make the	
	high-level recovery decisions.	158
5.13	Chapter 5 Highlights.	159
6.1	Illustration of Structural Differences (No substrate): (a) Planar Device; (b)	
	FinFET Device.	165
6.2	Cross section view of structural differences between (a) Bulk FinFET and	
	(b) SOI FinFET	166
6.3	Capacitance components for a FinFET device: (a) Cross-section view and	
	(b) Top view	169
6.4	Leakage current evolution with technology scaling	171
6.5	$I_{\rm on}/I_{\rm off}$ ratio with technology scaling	172
6.6	VTC curves under different supply voltages for a 1xnm FinFET inverter	
	(PMOS and NMOS are sized equally)	174
6.7	Velocity saturation index (α) for different technologies	176
6.8	Layout decomposition: A single layer is decomposed in two or more masks	
	to enhance the resolution.	177

6.9	Left side: a fin with a vertical slope which presents better short channel	
	metrics [129]. Middle: a standard fin with some degree of inclination as the	
	one used in the 22nm Intel's node [137]. Right side: a fin with a triangular	
	cross-sectional area that can help to reduce the leakage [65].	179
6.10	Different FinFET logic styles: 2-input NAND gate designs with SG and IG	
	devices	180
6.11	ON current vs. Body bias for a 2-finger 1xnm and 7nm NMOS transistor	181
6.12	16-input AND gate implemented with different stack height (1, 4 and 16).	182
6.13	16-input AND delay simulations with different stack height (interconnect	
	capacitance is considered).	183
6.14	Simulated FO4 delays for Inverter, 2-input NAND and 2-input NOR gates in	
	different technology nodes (all values are normalized to the 7nm FO4 INV	
	delay)	184
6.15	Simulated thermal characteristics (Delay vs. Temperature) in multiple tech-	
	nology nodes for a 9-stage ring oscillator. Blue - Super-threshold; Orange -	
	Near-threshold; Red - Sub-threshold	188
6.16	(a) Delay vs. V_{dd} ; (b) Energy/cycle vs. V_{dd} ; (c) Energy Delay Product (EDP)	
	vs. V_{dd} and (d) Minimum EDP values across multiple technology nodes	
	(simulated with the same NAND-based ring oscillator structure)	192
6.17	Transistor Aging: HCI occurs mainly during switching; PBTI happens when	
	NMOS is under stress; NBTI happens when PMOS is under stress. BTI	
	aging partially recovers during OFF states	197
6.18	Single transistor ON current degradation due to aging under DC nominal	
	voltage stress.	201
6.19	Single transistor ON current degradation due to aging under AC nominal	
	voltage stress with 50% duty cycle. Degradation is about half of the DC	
	stress case due to recovery	202
6.20	Simulation setup (an example of datapath): Aging can lead to timing failures	
	such as setup violation by slowing down the datapath. Designers should take	
	extra margins based on aging impact	203
6.21	Normalized Timing Margin vs. Temperature and Active Time: Margin shown	
	on Y-axis is normalized to the required aging margin for datapath (shown in	
	Figure 6.20) for 2 years at room temp (27° C)	204
6.22	Estimated aging margin for different IoT applications: X-axis corresponds	
	to IoT application index in Table 6.3, Y-axis shows the normalized design	
	margin and the error bars show design margin range within each category	205

6.23	6T SRAM read current degradation with aging for different temperatures	
	(nominal voltage)	206
6.24	IoT lifetime: Chip lifetime and Battery lifetime depend on different factors,	
	but they can affect each other indirectly. Two lifetimes together determine	
	the lifetime target of an IoT application.	208
6.25	Power and Aging profile of a typical IoT node: This figure is for illustration	
	only, the height and width are conceptually marked. For aging profile, Y-axis	
	"Aging" corresponds to aging-induced metric change such as ΔV_{th} or timing	
	margin (reduced performance).	211
6.26	Conceptual illustration of dynamic margins to enable one chip across multiple	
	IoT applications (BOL - Beginning of lifetime, EOL - End of lifetime)	212
6.27	Chapter 6 Highlights.	213
7.1	Illustration of Accelerated Self-Healing as a new dimension for mitigating	
	wearout effects.	216
B .1	Updated Top-down Design Flow with BTI Sensor Insertion.	226
B.2	Demonstration of the Sensor Insertion Flow in a Johnson Counter Design	227
B.3	The layout of the counter design after sensor insertion.	228

List of Tables

2.1	Summary of Test cases for Accelerated Wearout and Self-Healing	31
2.2	Summary of delay increase (%) under different temperatures	36
2.3	Parameter Descriptions for the Model	36
2.4	Summary of the Accelerated Self-healing results for 6 hours of recovery (%:	
	recovered percentage)	40
2.5	Summary of periodic accelerated rejuvenation test cases	49
3.1	Test Cases for EM Stress and Recovery	70
3.2	EM Line Current Limits @125°C for 10 years of operations	77
3.3	Temperature Derating Factor $\gamma_{der}(T)$ and EM Line Current Limit	77
3.4	Summary of Parameters for Black's Equation	77
3.5	Estimated EM stress and recovery time under normal operating conditions	
	for 28nm FDSOI technology	78
3.6	EM vs. BTI - Similarities and Differences during Stress and Recovery	82
4.1	Comparisons against other BTI sensor designs (all are circuit-level BTI sensors)	127
4.2	Summary of PPA Metrics for different Circuit Components (in 28nm FD-SOI	
	technology unless specified)	134
5.1	Simulated System Parameters	145
6.1	Summary of device parameters across multiple technology nodes (extracted	
	from I-V curves)	168
6.2	Normalized logical effort g and parasitic delay p values	184
6.3	Summary of IoT Applications Specifications (Aging-related Metrics)	199

Chapter 1

Introduction

1.1 Motivation and Background

The never ending demands for delivering higher performance, better energy efficiency and more integrations on a single die have enabled the continuous downscaling of CMOS transistor feature size. Although the technology scaling has been advantageous for many metrics, these advancements have also augmented the impact of the reliability issues [188]. Reliability refers to the probability that a system is able to perform its intended functions for a given lifetime under given conditions. Sources of unreliability at the hardware level include variations caused by manufacturing and operating conditions, soft errors caused by electrical noise or external radiations, and wearout (aging)¹ failures that are caused by device degradations [191]. Process variations have been well studied since many decades ago and are usually modeled accurately as part of the design kit, which guides circuit designers to design against it. Soft errors, also called transient faults or single-event upsets, may cause computation errors and corrupted data, but they are temporary and do not affect the lifetime of the computing systems [18]. Unlike the first two unreliability sources, wearout effects,

¹In this thesis, "aging" and "wearout" are used interchangeably. Circuit Aging/Wearout refers to wearout effects at both transistor level and interconnect level. Transistor Aging/Wearout mainly refers to BTI effect or other effects occur to CMOS transistors.

started to gain increasing attentions recently, are manifest during lifetime and highly depends on the unpredictable operating conditions. It has grown to be a huge reliability threat to the lifetime of digital circuits and systems [173, 202, 6]. Reasons behind this have been manifold [154], but it can be summarized in the following two categories.



Fig. 1.1 Wearout vs. technology scaling projected by Intel [151]. Y-axis refers to wearout related metric, and it is normalized to the 32nm node.

The first one is from the technology aspect, Figure 1.1 shows the projected wearout acceleration across multiple technology nodes from Intel [151], as technology scaling is reaching the nanoscale regime [6, 135]. The transistors become more susceptible to voltage stress [169, 61, 37, 135, 169] due to the increased effective field resulted from the reduced equivalent oxide thickness (EOT) [202]. Similarly, the shrinking geometries of metal layers render higher current densities, and the tremendous number of transistors within a compact area has resulted in higher power densities as well. Together, these lead to increased on-chip temperatures which potentially accelerate the wearout effects [86]. Moreover, advanced technology such as FinFETs have given rise to several new wearout concerns due to new

effects such as self-heating [165, 169]. As technology scaling enables billions of transistors fit on one chip, the challenge is that failure rate of one single transistor is required to decrease so that the historical values of mean time to failure (MTTF) of the whole system to be maintained.

Besides the technology scaling factors, wearout issues also become more pronounced from application perspectives. In high-performance computing applications such as servers, the system utilization has been increasing significantly especially due to the advent of the cloud computing. The goal for cloud operations is to maximize utilization by balancing compute jobs across an entire data center. That can lead to that a system runs most of the lifetime (> 5 years) without stopping. The approach is energy-efficient, but it can result in very high accumulated wearout and reduction of the performance eventually [223]. Utilization trends are shifting inside the emerging embedded applications such as the Internet of Things (IoT) edge devices and automobiles, which will continue until fully autonomous vehicles replace human drivers. In wearables or medical devices, where circuits usually work in near/sub-threshold for ultra low power (ULP) operation, the sensitivity of transistor ON current to threshold voltages is much higher than in super-threshold regimes [167]. Also, demanded by marketing and applications, these edge devices usually have very strict resiliency requirements [6, 54] and require longer lifetime. For example, some biomedical applications will require a lifetime of more than 50 years for medical implants [61, 54]. Finally, wearout issues are not just about time, they are highly thermal dependent. Many of such systems need to operate in extreme environmental conditions, such as high temperatures (without cooling), which, unfortunately, further accelerate wearout [9].

Wearout is a time-dependent reliability mechanism that is caused by several physical mechanisms that conspire to worsen metrics across the system hierarchy [78, 37], with performance degradation or intrinsic faults at the circuit level [159, 223], errors at the architecture level [80] and failures at the system level [80, 221]. It occurs at all parts of



Fig. 1.2 Illustration of BTI and EM. The Scanning Electron Microscope (SEM) figure (from [210]) shows a cross-section of a fully processed microprocessor. BTI occurs to transistors, and EM happens in metal wires.

a silicon chip. As shown in Figure 1.2, in general, at the transistor level, also known as Front-end of line (FEOL), Bias Temperature Instability (BTI) is one of the most prominent wearout mechanisms [135, 37, 20, 80]. It is characterized by the increase of the absolute value of threshold voltage ($|V_{th}|$) and the reduction of the carrier mobility (μ). In metal layers, known as Back-end-of-line (BEOL), Electromigration (EM) is the dominant reliability threat that increases the wire resistance *R* over time (soft wearout), and ultimately can break the wire (hard failure). EM is especially critical for power delivery networks (PDN) in modern ICs [208, 86, 184, 173]. Both wearout effects happen due to stress caused by voltage or current, when the stress is removed, there are some levels of recovery, but usually at much lower rate than the wearout process. In the next section, we will review some of the state-of-art techniques for mitigating both wearout effects.

1.2 Wearout Mitigation Techniques



Fig. 1.3 Taxonomy of BTI and EM Mitigation Techniques. Note that "Recovery Boost" here refers to the existing recovery acceleration solutions which mainly focused on BTI wearout for SRAMs.

In the past decade, various techniques have been proposed to deal with both BTI and EM wearout issues from system level down to circuit level, and from design time (static) to run time (dynamic). Overall, these techniques can be categorized into four categories as listed in Figure 1.3 (modified based on [106, 80]), each category includes several techniques that are applied at different hierarchies of a system stack. The most common solution for wearout issues is to tolerate wearout and add margins (guardband) during design time (pre-fabrication). However, predicting the margin under dynamic workloads and changing operating conditions is very difficult and many times unfeasible. Therefore, worse-case margin is commonly employed. Such timing margin can be >20% for 10 years [78, 97], and it can be as high as 14.5% for voltage margin [238]. The added margins mean large timing

slacks and therefore wasteful energy consumption ($\sim 30\%$ [106]) especially during the initial lifetime. The significance of these overheads also increases with scaling of technology nodes as discussed in the last section. A better circuit-design approach for tolerating wearout is by optimal sizing. Specifically, for BTI, upsizing the transistors compensates the V_{th} increase. EM effects are mainly addressed by design rules (e.g. increasing the metal width) during the physical design phase. While sizing is a very complex problem for solvers as it associates with multiple metrics, optimizing one can hurt another. Besides the area overhead, the increased transistor size contribute to the increased gate capacitance thus the dynamic power consumption, and it also increases the leakage power. Similarly, increased metal width contributes the load capacitance, which can increase the power consumption as well. Another design time method to address BTI degradation problems is during the technology mapping of the logic synthesis. The idea is to balance the delays of system components (gates, paths or even processor pipeline stages) by considering wearout so that the overall lifetime of the system can be optimized [107]. As this solution is highly based on prediction of wearout under dynamic conditions, it can lead to over-estimation and inaccurate results in many cases, which in turn leads to a low efficiency for design-time approaches. An alternative design time solution has been adding redundant resources for wearout-critical components such as critical path at the circuit level [15] or cores/processing units at the architecture/system level [200]. The overall lifetime can be improved by switching among redundant sources. While adding redundant elements can increase the area significantly, and can also lead to performance overhead during switching. It complexes the design process as well. Power gating has been a popular low-power design techniques that was originally used to lower the leakage power. It has been also adapted to help relieve BTI wearout as stress time is reduced while recovery time is increased [105]. But power gating only enables the passive recovery, which has been demonstrated being very slow and unpredictable (this will also be shown in Chapter 2 by our experiments). Compared to the BTI mitigation techniques during design

time, EM wearout has been dealt with by a less diverse group of techniques. The reason is that front-end designers (who design the architectures and IPs) have much less control over backend implementations, so physical designers usually need to take care of EM wearout by either upsizing the metal wire based on the design rule specified by the foundry, or by adding more metal straps to compensate the resistance increase caused by EM. This solution, although being used for many years, can lead to wasted routing resources and conservative design.

Compared to the static design time solutions, adaptive post-silicon techniques appear to be more "economic" in terms of costs and margins by compensating wearout during run-time. Previous work have proposed novel circuit and architecture level BTI [180] and EM sensors [77] to track and monitor wearout, and then several knobs can be adjusted correspondingly. Such knobs can be clock frequency [153], supply voltage [238], body bias [166] or combined [141]. At the system level, BTI-aware scheduling was proposed to equalize the utilization of functional units in a microprocessor to improve its lifetime reliability [192]. Although the dynamic margins enabled by these solutions can guarantee that the circuit and system is functioning in the presence of wearout, wearout itself is still unchecked, and the system can function but might run sluggish or burn more power gradually. Also due to the unique time dependent nature of the wearout (especially the irreversible components), which will accumulate fundamentally as the system runs. In many cases, wearout sensors (the expected number for a future chip can reach as many as hundreds [180]) need to be ON for tracking over the entire lifetime, and this will add unacceptable tracking power overhead. A better solution would be to somehow reduce the actual wearout induced variations by "repairing" them. Since BTI/EM are voltage/current dependent [217], one way is to reduce the voltage/current stress, thus to alleviate wearout during run time [238], but this way will introduce big performance overhead. The second way is to take advantage of the recovery properties of BTI/EM by generating more idle time for passive recovery (system unstressed when not in use) [74].

While passive recovery is very slow and unpredictable, and cannot be used to reduce margins effectively, thus it is sometimes even ignored when modeling wearout phenomena. Thus a solution that can fundamentally fix wearout instead of compensating for its effects would be clearly preferable. The concept of recovery boost was firstly introduced in SRAMs for cache blocks to recover some of the NBTI effects through reverse bias [191, 192], but these ideas were mostly on the intuition level due to lack of good understandings of device level recovery properties, also the implementation in these work can't be applied to other logics and didn't consider other wearout mechanisms such as EM. Different from all of these previous wearout mitigation solutions, in this thesis, we propose and demonstrate that both BTI and EM wearout can be made active by reversing the directions of the stress (e.g. using positive V_{gs} instead of negative V_{gs} for NBTI, using reverse current instead of forward current for EM) and can be accelerated (e.g. by increasing the temperature). Based on the actual hardware measurement results, these accelerated self-healing² techniques will lead to significant recovery rate improvements. On top of these, we investigate the irreversible components of both wearout effects, and propose a set of solutions which can completely mitigate and avoid wearout. On-chip implementations across the full system stack (circuit, architecture and system) to enable and assist both BTI and EM accelerated self-healing are also presented in details in this thesis.

1.3 Thesis Contributions

This thesis identifies the problem that future electronic systems will suffer wearout issues from both technology and application aspects, and there is no wearout mitigation solutions that can serve as a "panacea". The goal of this thesis is to explore an effective (in terms of power, performance and area) while being orthogonal (that can be used together with other solutions) dimension that is able to fundamentally repair BTI and EM wearout issues through

²In this thesis, "accelerated self-healing" and "accelerated and active recovery" are used interchangeably.

accelerated and active recovery. The main contributions of this thesis are summarized as follows:

1) Providing new and deeper understandings of BTI and EM recovery behaviors through comprehensive studies based on hardware experiments. Most of the wearout experimental research in the past looked into mainly the stress phase behaviors, recovery for both EM and BTI haven't been well understood. In particular, there are irreversible components for both wearout effects, and explanation of this has been manifold due to lacking of experimental studies. In this thesis, we study the recovery behaviors based on measurement on actual hardwares (FPGA for BTI, test chip for EM). Each set of measurement has been designed by considering different combinations of recovery conditions and lasts for more than 3 days. We also explore the boundary between the reversible and irreversible component of BTI and EM. The experimental results provide many new insights on recovery, such as frequency dependency, accelerated and active recovery and long-term recovery vs. short recovery difference. These insights can contribute as experimental evidences for reliability community to create better and more accurate device models for both wearout effects.

2) Demonstrating that accelerated self-healing is an effective solution for fixing wearout. Accelerated self-healing is essentially a "reverse" wearout process where several knobs are tuned during recovery to assist the process. We demonstrate that accelerated self-healing can lead to significant recovery rate improvement (e.g. 72.4% of the wearout is recovered within only 1/4 of the stress time for BTI; 70% of EM wearout can be recovered within 1/5 of the stress time). Our further study demonstrates that there is still a lingering permanent component are irreversible even under accelerated and active recovery conditions, we explore the boundary between reversible and irreversible wearout physically and experimentally. By studying the frequency dependent behaviors of BTI and EM wearout and recovery, we demonstrate that the boundary is actually "soft" and can even be controllable,

and this leads to a biology-inspired sleep-when-getting-tired strategy that keeps the circuit active only during the reversible phase of wearout until the irreversible wearout kicks in, thus the irreversible wearout becomes almost unobservable even in accelerated stress cases. The proposed accelerate self-healing cases, together with scheduled explicit accelerated self-healing periods ahead of any sign of stress, will be simpler to implement on chip and results in that the system operating for a longer time in a "refreshed" mode, thus leading to better performance, and has better cumulative metrics (e.g. average performance) as well.

3) Implementing accelerated self-healing on-chip across the full system stack. To enable the implementation of accelerated self-healing techniques and fully utilize the unique recovery behaviors for BTI and EM wearout, we propose a full set of potential implementation solutions at the circuit, architecture and system levels. These solutions cover all aspects from recovery-driven design methodology, novel sensing techniques for monitoring both wearout and recovery, novel power gating structures to recovery assist circuit for enabling multiple recovery modes. Some of these circuit solutions have also been successfully demonstrated on test chips. Since single-layer recovery solution is not the most cost-friendly, we also discuss the implications of implementing accelerated self-healing at different levels of a system hierarchy. Novel schemes and scheduling solutions have been presented. The recovery circuit components proposed in this thesis are flow friendly, and can serve as the infrastructures for future research in this direction to build upon. Combining all these techniques enables an true accelerated self-healing system that benefits from the full-recovery capabilities.

4) Providing a comprehensive technology study across multiple nodes, and studying FinFET wearout in IoT applications. As FinFET has become a main-stream technology node for most of the high-performance computing chips, it is also in the process of being adapted to low-cost embedded systems, in this thesis, we perform a comprehensive technology study across multiple nodes ranging from planar, FD-SOI to FinFET based on both foundry provided models and predictive models. The study explores new design challenges and new
insights since the advent of FinFET technology. It can be contributed as an educational material and design guide for circuit designers who design with FinFETs. At advanced technology node, wearout has become more pronounced, and this can especially be critical in IoT context in which some applications require very reliable operations (zero-error tolerance) spanning a much longer lifetime (> 50 *years*). We look into this aspect as well in this thesis by studying how transistor wearout can impact different categories of IoT applications with the foundry-provided wearout models. We conclude that wearout needs to be considered in the full design cycle and the IoT lifetime estimation needs to incorporate wearout as an important factor. We also present application-specific solutions to mitigate wearout in IoT systems.

1.4 Thesis Organization

The remainder of this thesis is organized as follows. The thesis organization is also given in Figure 1.4.

Chapter 2 presents the BTI accelerated self-healing measurement results on the FPGA platform. Details of experimental setup, test cases and test results are covered in this chapter. An analytic gate-level BTI stress and recovery model is described. Based on the measurement and the model, we show the impacts of applying accelerated self-healing techniques on other digital systems.

Chapter 3 presents the EM accelerated and active recovery measurement results with on-chip metal lines. We describe the test flow, test structure and experimental results, which also lead to a thorough analysis on EM signoff when applying the recovery techniques. EM recovery and BTI recovery behaviors are also compared in this chapter.

Chapter 4 describes a set of circuit IPs for implementing the idea of accelerated selfhealing on chip. These IPs include on-chip negative voltage generators, on-chip heat generators, multi-mode EM/BTI recovery scheme and three types of EM/BTI sensors. We



Fig. 1.4 Thesis Organization.

present the simulation results and physical implementations for each circuit and analyze the overheads of the designs.

Chapter 5 explores the cross-layer implementations of accelerated self-healing. We discuss how to select the recovery knobs in real systems and how to instrument recovery at the architecture level and system level. Various potential implementations are described. We also comment on potential overheads of each implementation and how cross-layer techniques can work together to leverage the overheads.

Chapter 6 can be divided into two parts. In the first part, we study the circuit metrics across multiple technology nodes, including planar, FDSOI and FinFET. The study shows that FinFET technology is superior in many dimensions, but there are also newly introduced challenges to be addressed during design. The second part presents a study on impact of

FinFET transistor wearout in IoT domains. We classify current IoT applications based on wearout-related metrics and conduct circuit wearout simulations for each category. Potential solutions for mitigating the effect of wearout in IoT circuits are also presented in this chapter.

Chapter 7 concludes the thesis. We summarize the major contributions and also discuss the future directions in this chapter.

Chapter 2

Accelerated Self-Healing Techniques for BTI Wearout

2.1 Overview

Bias temperature instability (BTI) is one of the most dominant wearout mechanisms for transistors, especially in advanced technology nodes [134]. It increases the threshold voltage (V_{th}) and reduces the mobility (μ) of transistors over time under voltage stress, thus increasing the circuit delay and necessary time margins [135, 223]. As shown in Figure 2.1, Negative Bias Temperature Instability (NBTI) occurs under negative stress conditions and affects PMOS transistors. Similarly, Positive Bias Temperature Instability (PBTI) affects NMOS transistors under positive stress voltage. Depending on the bias condition of the gate, there are two phases of BTI. The stress (or wearout) phase is defined when the gate is under voltage stress ($V_{gs} < 0$ for PMOS, $V_{gs} > 0$ for NMOS), and the *passive* recovery phase happens when transistors are in OFF state, where voltage stress is "released" ($V_{gs} = 0$). While passive recovery has been accepted as being slow and unpredictable, and cannot be used to reduce margins effectively, thus it is sometimes even ignored when modeling wearout phenomena and estimating the guardband. Different from previous solutions, in this chapter,



Fig. 2.1 BTI Stress and Passive Recovery: NBTI happens in PMOS transistors; PBTI happens in NMOS transistors. For both mechanisms, BTI starts recovering when transistors are turned OFF, but this passive recovery period is very slow and unpredictable.

we demonstrate that BTI recovery can be made active by reversing the direction of the voltage stress (e.g. using positive instead of negative V_{gs} in the case of NBTI), and can be accelerated (e.g. by increasing the temperature), thus leading to accelerated self-healing. Based on the actual hardware measurement results on 40nm FPGAs, these accelerated self-healing techniques will lead to significant recovery rate improvement (e.g. we demonstrated a case where 72.4% of the wearout is recovered within only 1/4 of the stress time). We also explore that even in the accelerated active recovery case, there are still components of BTI wearout that are irreversible. By studying the frequency dependent behaviors of BTI wearout and recovery, we demonstrate that the boundary is actually "soft" and can even be controllable, and this leads to a biology-inspired sleep-when-getting-tired strategy that keeps the circuit active only during the reversible phase of wearout until the irreversible wearout kicks in, thus the irreversible wearout becomes almost unobservable even in accelerated stress cases. The proposed solutions, with proactively scheduled explicit accelerated self-healing periods ahead of any sign of stress, will be simpler to implement on chip and results in that the system operating for a longer time in a "refreshed" mode, thus leading to less necessary BTI-induced margin, better performance, and has better cumulative metrics (e.g. average performance) as well. The details of experimental setup, measurement results and modeling of the proposed accelerated self-healing solutions are presented in this chapter.

2.2 BTI Wearout and Recovery Basics

The mechanisms behind BTI have been quite controversial and still haven't reach a consensus [203, 136]. There are mainly two types theory to explain BTI effects, and they are illustrated in Figure 2.2. The first one is called "Reaction-Diffusion (RD)" theory and is shown in Figure 2.2a [216]. According to this theory, BTI has been attributed mainly to interface trap generation. For example, for NBTI case (PMOS), when the transistor is ON, the voltage stress across gate and source could potentially break the covalent bound of $S_i - H$ at interface, this process is called reaction. The separated hydrogen atoms combine to form H_2 , which diffuses toward the gate of the transistor. These broken $S_i - H$ bonds generate positively charged traps for holes and leads to transistor parametric shift, such as increased threshold voltage V_{th} . When the PMOS switches off (i.e., $V_{gs} = 0$), and stress is removed, the recovery phase starts, where holes are not present in the channel and thus, no new interface traps are generated; instead, H diffuses back and anneals the broken Si-H. As a result, the number of interface traps is reduced during this phase and the NBTI degradation is recovered passively.



(a) BTI Reaction-Diffusion (RD) Mechanism(b) BTI Trapping-Detrapping (TD) MechanismFig. 2.2 Two BTI Mechanisms.

While the R-D based theory can well predict the constant stress behaviors, recent advances in fast on-chip BTI measurement have explored several dynamic BTI behaviors that can't be explained and are inconsistent with what has been modeled by R-D theory [218]. This led to an alternative way to explain the BTI effect with trapping/detrapping (TD) mechanism, and it has been validated against silicon and become widely used in the community [37, 87, 135, 143]. As shown in Figure 2.2b, when the PMOS transistor is ON, the trap energy (relative to Fermi energy level) is modulated. If the trap gains sufficient energy, it may capture a charge carrier, thus reducing the number of available carriers in the channel, the charged trap state modulates the V_{th} and acts as a scattering source, reducing the effective mobility. This is called Trapping process. Similar to what has been captured in RD theory, if the transistor is OFF and in passive recovery phase, some of the interface traps can be annealed slowly (shown as De-Trapping process in the figure), and the number of occupied traps reaches a new equilibrium and results in partial recovery. Although the effect of PBTI has been negligible in previous technologies, it is rapidly becoming an important reliability issue with the introduction of high-k and metal gates [19, 157, 237]. Since the degradation effect of PBTI is similar to NBTI, the PBTI effect can be modeled similar to the NBTI effect [237].

There have been also recent modeling frameworks which incorporate both RD and TD theory [69, 158], but these work show that the most dominant component of BTI is due to hole trapping and interface trap generation. Thus the presented modeling work is mainly based on the Trapping/Detrapping theory in this chapter. The device-level TD model developed in [217] captures the details of the stress and passive recovery physically. According to the model, the threshold voltage of the transistor increases logarithmically, and the overall dynamic BTI behaviors can follow is shown in Figure 2.3. Assume that the single transistor is turned on (stress period starts) at time t = 0, and no voltage stress is applied before. The threshold voltage increase (ΔV_{th}) until time t_{st} is modeled as:

$$\Delta V_{th}(t_{st}) = \phi_{st} \left(A + log(1 + Ct_{st}) \right)$$
(2.1)



Fig. 2.3 BTI behaviors modeled by Trapping/Detrapping (TD) theory. Passive Recovery still leaves a net ΔV_{th} that is almost permanent and hardly recovered within a reasonable time.

If a recovery interval of t_{rec} follows the stress phase, the total threshold voltage shift in the end is equal to:

$$\Delta V_{th}(t_{st} + t_{rec}) = \phi_{rec} \left(A + \log(1 + Ct_{rec}) \right) + \Delta V_{th}(t_{st}) \left(1 - \frac{A + \log(1 + Ct_{rec})}{A + \log(1 + C(t_{st} + t_{rec}))} \right) \quad (2.2)$$

$$\phi \sim Kexp(\frac{-E}{kT})exp(\frac{BV_{dds/r}}{kTt_{ox}})$$
(2.3)

where A, B, C are (approximately) constant across the same technology node, K is the fitting parameter, k is Boltzmann's constant, T is temperature, E is activation energy, t_{ox} is the oxide thickness, and V_{dds} and V_{ddr} are the supply voltages under stress and recovery, respectively. More details of the model can be found in [217].

The device-level BTI model discussed above indicates the strong dependence (exponentially) of the threshold voltage shift on voltage and temperature during both stress and recovery period. It serves as the basis of our circuit and higher level accelerated self-healing model framework. All parameters and their values used in the model are extracted based on our measurements which will be discussed in details in the following sections.

2.3 Prior Work on BTI Recovery

The partial recovery property of BTI has been utilized in many work to improve the lifetime and other metrics (e.g. performance) of the digital systems. Several methods [3, 43, 73, 74, 67] were proposed to rebalance the signal probabilities for logic or SRAMs to maximize the passive recovery time at the circuit level. An alternative method was to adaptively tune the performance according to the degree of BTI wearout so that certain blocks could start the passive recovery phase earlier [212, 189]. Novel schemes were proposed to exploit the idle time of busy functional units for out-of-order processors [192] and superscalar architectures [124]. A dynamic routing algorithm was proposed to adapt to the wearoutinduced degradations in heterogeneous NoCs [13]. Since passive recovery is much slower than the wearout process [172], wearout gating [28] was firstly proposed, the idea was to add a coupling transistor to a regular power gating structure so that the virtual supply and the logic high/low supply are equalized, in this way, the voltage across the logic block is zero and transistors experience zero stress. The goal of this method was still to release the stress completely for passive recovery. To further boost the recovery process, intense recovery for logics [29] and SRAM array [191, 193] was then proposed, the idea was to raise the gate voltages of a chain of logic or memory cell in order to put PMOS devices into the recovery enhancement mode in which the voltage across the transistor is reverse bias mode with full range of V_{dd} . This method can potentially be harmful to device reliability since it can lead to device breakdown under a high reverse voltage, it also incurs very high power routing costs and design complexity. A power napping scheme was proposed to help the recovery of NBTI and PBTI [16]. But all of these previous work focused on only SRAM cell designs or architectural level implementations, it was still unclear how much benefit recovery boost



Fig. 2.4 Proposed accelerated self-healing solutions for NBTI. Similar solutions can be applied to PBTI as well.

could achieve due to lack of experimental data and models. In addition, these solutions still leave the irreversible wearout unchecked. Wafer level and transistor level experiments and theory [5, 99, 162], together with the device level model discussed in above section, indicated that BTI recovery highly depends on temperature, these work provided a physical evidence for our proposed accelerated self-healing techniques. Several recent work [70, 52, 158] have studied the irreversible component of wearout at the device level. However, these work focused only on demonstrating and modeling the permanent component, thus a solution that could fundamentally repair the irreversible wearout is still missing in the field. The solutions presented in this chapter differ from the previous recovery and recovery boosting work in several aspects. Firstly, we demonstrate that both high temperature and negative voltage could accelerate recovery with actual hardware measurements on 40nm FPGAs and also develop the analytic model for them. Secondly, we propose a biology-inspired strategy that runs the system in such a way that the irreversible wearout can be fully avoided so that the system keeps almost "fresh" all the time. Lastly, the proposed solutions can be effectively applied to all logic blocks, instead of just SRAM arrays. The notion of accelerated self-healing can be implemented across the system hierarchy and can be introduced as a key design knob for cross-layer resilience.

2.4 Accelerated Self-Healing

2.4.1 Active and Accelerate BTI Recovery

In this thesis, we postulate that the systems will use sleep time as an active recovery period essential for their overall performance. We demonstrate that during sleep, several *accelerated self-healing* solutions can be implemented to deeply rejuvenate the circuit, and they are shown in Figure 2.4 (for NBTI recovery as an example). Firstly, BTI recovery can be made *active* by "turning off" the transistor more via a negative voltage across the source and gate. Secondly, high temperature can increase the kinetic energy for the charge carriers, thus leading to the *accelerated recovery*. The third case is when active recovery can be further accelerated through the joint effort of both negative voltage and high temperature. From physics perspectives, the self-healing techniques reverse the direction of BTI wearout and increase the recovery rate. The applied negative voltage (active recovery) across the transistors activates the detrapping process by pushing trapped carriers back to their original states, while high temperature can further assist the healing process by exciting the carriers. Overall, two solutions can work together to achieve the highest possible recovery rate.

2.4.2 Gate-level Analytical Model for Accelerated Self-Healing

To model the performance degradation and rejuvenation due to wearout and accelerated self-healing, a gate level analytical model that is based on the device model describe in (2.1)-(2.3) is developed in this thesis. The model serves as a connection between circuit metrics (such as delay t_d) and device level parameter shift (mainly V_{th}) shift. Based on the circuit theory [167], the propagation delay of a digital gate can be approximated as:

$$t_d \sim \frac{C_L V_{dd}}{I_d} \propto \frac{C_L V_{dd}}{V_{dd} - V_{th}} \tag{2.4}$$

where C_L is the output capacitance of the gate. The change in gate delay when V_{th} is subject to change is:

$$\Delta t_d \sim \frac{\Delta V_{th}}{V_{dd} - V_{th}} \cdot t_{d0} \tag{2.5}$$

where t_{d0} is the original delay of the gate without any V_{th} shift. Combine equations (2.1), (2.3) and (2.5), the total delay increase after t_{st} can be expressed as:

$$\Delta t_d(t_{st}) = \beta_{st} exp(\frac{-E}{kT}) exp(\frac{BV_{dds}}{kTt_{ox}}) \left(A + log(1 + Ct_{st})\right)$$
(2.6)

where β_{st} , *A*, *B* and *C* are fitting parameters and can be extracted from measurement results. During the accelerated self-healing phase, based on the recovery phase equation of the device model, we combine equations (2.2), (2.3) and (2.5), and the delay change after sleep period t_{rec} becomes:

$$\Delta t_d(t_{st} + t_{rec}) = \frac{\phi_{rec}}{V_{dds}} \left(A + \log(1 + Ct_{rec}) \right) + \Delta t_d(t_{st}) \left(1 - \frac{A + \log(1 + Ct_{rec})}{A + \log(1 + C(t_{st} + t_{rec}))} \right) \quad (2.7)$$

Assume that the ratio of operation time to active sleep time of the system is α , the total time (stress time + sleep time) is t_{total} , based on (2.7), the overall delay increase will be:

$$\Delta t_d(t_{total}) = \phi_{acce} \left(A + log(1 + C\frac{t_{total}}{1 + \alpha}) \right) + \Delta t_d(\frac{\alpha t_{total}}{1 + \alpha}) \left(1 - \frac{A + log(1 + C\frac{t_{total}}{1 + \alpha})}{A + log(1 + Ct_{total})} \right)$$
(2.8)

$$\phi_{acce} \sim \frac{K}{V_{ddr}} exp(\frac{-E}{kT_{acce}}) exp(\frac{BV_{ddr}}{kT_{acce}t_{ox}})$$
(2.9)

Several observations can be made based on the gate-level accelerated self-healing model. Firstly, the exponential dependency of the total delay increase on recovery temperature and voltage shows that by increasing T_{acc} and decreasing V_{ddr} , the first component in Equation (2.8) can decrease significantly. The second observation is that the active vs. sleep ratio α also affects the overall delay change. The final observation is that the recovered part $(\Delta t_d(t_{total}) - \Delta t_d(t_{str}))$ highly depends on the previous stress history $(\Delta t_d(t_{str}))$.



Fig. 2.5 Illustration of a potential use case flow for Gate-level Accelerated Self-healing model. This flow can be used for evaluating how recovery conditions can affect the circuit metrics and system lifetime further.

This model captures the circuit-level metric change due to wearout and accelerated and active recovery. It can be potentially integrated into higher-level abstract models for exploring the accelerated self-healing space, and this is illustrated in Figure 2.5. The delay change due to wearout can lead to timing violations at the circuit level, and this could lead to failures at the system level. Based on the recovery conditions, the gate-level model is able to estimate the "recovered delay" which can potentially be used to predict the reduction of failure rate and extension of the system lifetime. Several use cases of this model will be detailed in the following sections.

2.5 Experimental Setup

2.5.1 Test Platform

FPGA vendors have been aggressive in adopting the very latest technology nodes; this makes FPGAs more susceptible to wearout that can lead to frequency degradation [168]. Due to their "bleeding edge" technologies, reconfigurability and regular structure, FPGAs are an ideal test platform for wearout research [182, 112]. In this thesis, we choose 2-input Look Up Table (LUT)-based commercial FPGAs [89] fabricated in a 40nm technology to



Fig. 2.6 FPGA test platform (Lattice Semiconductor iCE40 HX-Series) architecture [89].

demonstrate experimentally the proposed accelerated self-healing techniques. Figure 2.6 shows the architecture of the FPGA. Basic components of FPGAs include the I/O and the core architecture; we use only the core architecture for testing wearout in this thesis. The core includes 1280 logic cells (LUT + Flip-flops), which are grouped in Programmable Logic Blocks (PLB) which can be programmed to perform logic and arithmetic functions. Each PLB consists of 8 interconnected logic cells as shown in the figure. The chip has a SPI port that supports programming and configuration of the device using the standard FPGA synthesize and Place & Route (PnR) flow.



(a) BTI Test Structure: A 75-stage Ring Oscillator mapping to the Look-Up-Table (LUT) structure on FPGAs.



(b) A floorplan showing how BTI test structure is mapped on FPGA fabric: green squares represent utilized logics, red represents interconnect, I/Os interface with the host mother board and programmer through flat cable.

Fig. 2.7 BTI Test Configuration on FPGAs.

2.5.2 Test Configuration

In the gate-level model discussed in Section 2.4.2, the delay change is used as the metric to capture the effect of wearout and recovery, and the same metric is also employed in the experimental part. We choose a Ring oscillator (RO) structure which is widely used as a test platform [149] to measure the delay of the Circuit Under Test (CUT) to capture the delay change. Figure 2.7a shows the test configuration, which is a modified LUT-based Ring Oscillator based on the design proposed in [219]. It consists of 75 inverters implemented

in LUTs and a 16-bit counter to capture the output frequency of the ring oscillator. Enable signal En is used to switch between AC stress (switching) and DC (constant) stress mode. Figure 2.7b shows how the test structure is mapped to the FPGA fabric. The CUT is kept under voltage stress and it is enabled only every 20 minutes for frequency recording. The oscillation frequency f_{osc} can be calculated as:

$$t_d = \frac{1}{2f_{osc}} = \frac{1}{4f_{ref}C_{out}}$$
(2.10)

where f_{ref} is the frequency of the reference clock. To pick this frequency, CUT is placed at different locations on the FPGA as shown as different indexes in Figure 2.7b (Ring X), and a diagnostic program is run. The output of the counter is read from a certain time range that has stable values. Environmental factors and the voltage supply are kept constant from one reading to another; when $f_{ref} = 500Hz$, the variation of the counter output is within ± 5 and $\pm 0.0001\%$ in terms of the corresponding RO frequency variation which we consider acceptable.

2.5.3 Test Flow and Test Conditions

Test Flow

To sample the stress and recovery data from the FPGA chip in a fast and efficient manner, an automatic test flow is developed and is shown in Figure 2.8. The FPGA chip interfaces with the PC through an evaluation board (AT91SAM7SE-EK) developed by Atmel [139]. The micro-controller unit (MCU) on the board can be programmed in C and serves as a "signal generator" for the control signals in the BTI test structures on FPGA chip. The timer/counter in MCU is employed to control how long the FPGA chip is stressed and recovered, it also controls the analog switch which can pass or cut the supply voltage for the test chip. On the FPGA side, the RTL description of the BTI test structure (shown in Figure 2.7a) is loaded

into FPGA through a programmer and Serial Peripheral Interface (SPI). The reference clock for the counter is given by an external clock generator. Counter outputs are read out and saved to PC through the mother board interface. By using this test methodology, the total data sampling overhead is less than 3s, which is believed to have a negligible impact on overall BTI wearout and recovery behaviors.



Fig. 2.8 BTI Test Flow: FPGA chip communicates with computer through an Atmel evaluation board, on which a micro-controller can be programmed with C. The inherent timer in MCU can be programmed to controller how long the chip is stressed or recovered.

Stress and Recovery "Knobs"

In Section 2.4, we have introduced three ways to accelerate the recovery during sleep, one is active recovery through small negative voltage, another is the high temperature and the third is the combination of the first two. Thus the main "knobs" tuned in measurements are voltage, time, temperature, switching activity and α , the ratio of stress (active) and recovery

(sleep or rejuvenation) time. Two stress modes are considered – AC stress and DC stress (constant stress). DC stress refers to the case where the input of the test structure is always static and doesn't switch, and AC stress is when the input is switching with 50% duty cycle (*En* signal shown in Figure 2.7a can switch between modes).

Accelerated Test Methodology

Since both elevated temperature and voltage have a great impact on wearout and can be used to accelerate wearout. In our experiment, both elevated temperature and high voltages are applied so that we can observe a larger than 1% frequency degradation under high temperature for all test cases. The recommended operating temperature of the FPGAs we use is within -40° C to 85°C. In our test cases, we use 100°C and 110°C, which are above the upper limit of temperature, but not too high to prevent the chip from functioning. The FPGA chips are heated up or cooled down by a thermal chamber, which allows temperature fluctuation of $\pm 0.3^{\circ}$ C. Core voltage is provided by a DC power supply and its nominal value is 1.2V, the elevated voltage is within the 10% range and is much lower than the breakdown voltages. Figure 2.9 shows the measurement setup. FPGA chip is placed in a zif socket board which connects with the interface board through a flat cable physically to ensure that only the FPGA chip is exposed to high temperature environment.

Test Cases

All tests are carried out on a group of fresh commercial FPGA chips. Several test cases in both stress (wearout) and recovery (including accelerated self-healing) phases are considered and are denoted as follows (AS – accelerated stress, AR – accelerated recovery, etc.):

• AS110AC24: In this accelerated stress test case, the chip is under 110°C environment for 24 hours in AC stress mode. RO is always enabled to switch.



Fig. 2.9 BTI Test Setup: FPGA chip is placed in a socket board, which connects with the interface board through a flat cable to separate the micro-controller from the high temperature environment, only the FPGA chip is in thermal chamber.

- AS110DC24: This is similar to the previous case, but in DC stress mode. RO is enabled only every 20 minutes for data recording. Data sampling overhead is less than 3s.
- AS100DC24: 100°C is applied and the chip is under accelerated DC stress mode for 24 hours.
- 20Z6: In this case, chips are recovered for 6 hours under 20°C at 0V.

- **AR20N6**: Negative voltage of -0.3V is applied to the chip to activate the recovery at 20°C.
- **AR110Z6**: In this case, only high temperature (110°C) is applied, and the chip is powered off at 0V for 6 hours.
- **AR110N6**: Chips are recovered with both 110°Cand -0.3V.

Phase	Case No.	Chip	T (° C)	V _{dd} (V)	Time	Mode	Stress Time	
		No.	(-)	- uu (·)	(hours)		Recovery Time	
	AS110AC24	1	110	1.2	24	AC	-	
	AS110DC24	2	110	1.2	24	DC	-	
Stress	AS110DC24	3	110	1.2	24	DC	-	
(Active)	AS100DC24	4	100	1.2	24	DC	-	
	AS110DC24	5	110	1.2	24	DC	-	
	AS110DC48	5	110	1.2	48	DC	-	
Recovery (Sleep)	R20Z6	2	20	0	6	-	4	
	AR20N6	3	20	-0.3	6	-	4	
	AR110Z6	4	110	0	6	-	4	
	AR110N6	5	110	-0.3	6	-	4	
	AR110N12	5	110	-0.3	12	-	4	

Table 2.1 Summary of Test cases for Accelerated Wearout and Self-Healing

During recovery, a negative voltage and/or a high temperature of 110°C are applied. During both stress and recovery periods, the test structure is enabled from the stress phase every 30 minutes for data sampling. The chip that is stressed and recovered under normal conditions (T = 20°C, $V_{dds} = 1.2V$ and $V_{ddr} = 0V$) is used as the baseline for comparisons. All test cases are summarized in Table 2.1.

2.5.4 Modeling BTI Stress and Recovery for FPGA Test Structures

In Section 2.4.2, a gate level analytic model that converts the device level parametric shifts to circuit level metric changes (delay change) was discussed. To analytically understand how BTI accelerated stress and recovery can affect the specific test structure on FPGA we chose, we apply the gate level model to Look-up-Table (LUT) circuit, details are as follows.

As shown in Figure 2.6, the basic building block of the FPGA core is LUT, which is mapped as inverter logic in our BTI test structure (Figure 2.7). Shown in Figure 2.10 is a generic Pass Transistor (PT)-based 2-input LUT structure. Routing blocks include all the routing elements between LUT blocks. Four configure bits (C0 to C3) are stored in Block RAM (BRAM), In0 and In1 are input signals. Let's consider an inverter mapped to the LUT: In0 is the input of the inverter, C0 to C3 are 0101 and In1 is always 1. As shown in the figure, the Path Of Interest (POI) is from the input of the LUT-based inverter to the output of the routing blocks. Assume the inverter is under DC stress, and In0 is always 1. M1, M5 are under stress and the threshold shift will affect the delay of POI. If In0 is always 0, only M7 is under stress. Based on this simple example, two hypotheses can be made:

- Not all the transistors on POI are under stress. In DC stress mode, once the inputs are given, the number of stressed and unstressed transistors is constant;
- Recovery can only have an impact on stressed transistors, but has no effect on "fresh" (never aged) transistors, nor on transistors that have already recovered (close) to the "fresh" state.

Although the exact gate level netlists of commercial FPGAs are unavailable, we believe that the two hypotheses can be applied to any pass-transistor LUT configurations.

Assume that all stressed transistors on POI are under the same stress condition (V_{gs} are the same), so we can approximately assume that ΔV_{th} of all stressed transistors are the same. Total delay change ΔT_d of POI becomes:

$$\Delta T_d = \sum_{n}^{LD} \Delta t_{dn} \sim \Delta t_d N_s \tag{2.11}$$

where *LD* is the logic depth, N_s is number of transistors that are under stress and $0 \le N_s \le LD$, Δt_d is the delay change for a single gate from Equation 2.6. Combine Equation 2.6 and 2.11,



Fig. 2.10 Pass-transistor based LUT structure.

and assume $V_{dds} \gg V_{th}$, the total delay shift of a whole path after a stress period of t_{st} can be expressed as:

$$\Delta T_d(t_{st}) = Yexp(\frac{-E}{kT})exp(\frac{BV_{dds}}{kTt_{ox}})(A + log(1 + Ct_{st}))$$
(2.12)

$$Y \sim K_{st} N_s t_{d0} \tag{2.13}$$

If V_{dds} and T are constant over stress duration, Equation 2.12 can be expressed as:

$$\Delta T_d(t_{st}) \sim \beta (A + log(1 + Ct_{st})) \tag{2.14}$$

where β , *A* and *C* are fitting parameters and can be extracted from measurement results. Similarly, during accelerated recovery phase, we combine Equations 2.8, 2.9 and 2.11 and calculate the delay change of POI after recovery period t_{rec} as:

$$\Delta T_d(t_{st} + t_{rec}) = \frac{\phi_{rec}}{V_{dds}} \left(A + \log(1 + Ct_{rec}) \right) + \Delta T_d(t_{st}) \left(1 - \frac{A + \log(1 + Ct_{rec})}{A + \log(1 + C(t_{st} + t_{rec}))} \right)$$
(2.15)

Assume the the ratio of operation time to active sleep time of the system is α , the overall delay change in one cycle will be:

$$\Delta T_d(t_{total}) = \Phi_{acce} \left(A + log(1 + C\frac{t_{total}}{1 + \alpha}) \right) + \Delta T_d(\frac{\alpha t_{total}}{1 + \alpha}) \left(1 - \frac{A + log(1 + C\frac{t_{total}}{1 + \alpha})}{A + log(1 + Ct_{total})} \right)$$

$$(2.16)$$

$$\Phi_{acce} \sim K_{acce} exp(\frac{-E}{kT_{acce}}) exp(\frac{BV_{ddr}}{kT_{acce}t_{ox}})$$

$$(2.17)$$

The model evaluation and validation will be discussed together with the measurement results in the following sections.

2.6 Test Results for Accelerated BTI Wearout

This section presents the testing results during the stress period. It shows the performance degradation under various stress conditions.

2.6.1 AC Stress vs. DC Stress

AC stress and DC stress are conducted in the first and second case described in Section 2.5.3, Figure 2.11 shows the measurement results at 110°C. The AC stress with a 50% duty cycle can be treated as a symmetric stress vs. passive recovery case. In the first 3 hours, RO frequency degradation of both cases is relatively fast and then becomes slower. AC stress has a symmetric stress and recovery process, during which stress phases are always



Fig. 2.11 AC and DC Stress Measurement Results: AC stress case with a 50% duty cycle shows slower BTI wearout.

followed by recovery phases due to dynamic activity of the circuit, and results in smaller frequency degradation, which is about half of that in the DC stress case. The results indicate that passive recovery is much slower compared to wearout since the chip cannot be fully recovered with symmetric AC stress. In other words, AC stress is a only partially self-healing process with a very slow recovery rate. To almost fully rejuvenate the chip, accelerated self-healing techniques are required.

2.6.2 Effect of Temperature on BTI Wearout

Figure 2.12 shows measured delay change over time at 100°C and 110°C. As the model predicts, initially, frequency degrades fast and then slower. High temperature accelerates the degradation. Table 2.2 summarizes the delay change (%) for different temperature conditions. Table 2.3 shows the extracted parameters we use in the model that is discussed in Section 2.5.4.



Fig. 2.12 Accelerated BTI Wearout under 110°C and 100°C for 1 day

Tomporatura (°C)	Measu	rement	Model Prediction		
Temperature (C)	12 hours	24 hours	12 hours	24 hours	
20	0.13%	0.19%	0.11%	0.18%	
100	1.1%	1.5%	1.1%	1.53%	
110	1.45%	2.16%	1.57%	1.96%	

Table 2.2 Summary of delay increase (%) under different temperatures

Table 2.3 Parameter	Descriptions	for the Model
---------------------	--------------	---------------

Parameter	Description	Value
k	Boltzmann Constant	1.38×10^{-23} J/K
E	Activation Energy	0.49eV
t_{ox}	Oxide Thickness	1 <i>nm</i>
А	Constant	0.2801
В	Constant	3×10^{-29}
С	Constant	0.8614
K _{st}	Fitting Parameter	$4.7 imes 10^{-4}$
Kacce	Fitting Parameter	7.34×10^{-5}

2.7 Test Results for Accelerated Self-Healing Techniques

This section will present the test results for the proposed accelerated self-healing techniques. To make comparisons, we use recovered delay (RD – delay decrease during recovery) as our metric, which can be calculated as $RD(t_2) = T_d(t_1) - T_d(t_2) = \Delta T_d(t_1) - \Delta T_d(t_2)$, where $t_2 > t_1$, and ΔT_d is the delay change with respect to the delay at time zero.

2.7.1 Negative Voltage

Currently, when electronic systems go to sleep, the supply voltage is usually gated to reduce leakage, but this only results in *passive* recovery. For *active recovery* we apply a negative voltage. The challenges for picking this negative voltage are:

- Breakdown voltage limitation: the voltage must be at the level below the lateral pn-junction breakdown voltage;
- Implementation feasibility: implementation of negative voltage will introduce area overhead;
- Gate-induced Drain Leakage Current (GIDL) may introduce large leakage currents.

In this test, we picked a negative voltage of -0.3V that is validated to be still within the "safe" margin of the breakdown voltage and leakage. Figure 2.13 compares the recovered delay over 6 hours (1/4 of the total stress time) when the temperature is set at 20°C and 110°C, respectively. Model predictions from Section 2.4.2 are also included in the figure. The coefficients used in the model are all extracted from the measurements. The results in Figure 2.13 show that stressed chips rejuvenate faster with a negative supply voltage for both temperatures. By applying a negative voltage the recovery is significantly accelerated even at room temperature.

2.7.2 High Temperature

High temperature not only accelerates wearout, this section will show that it will also accelerates recovery. Figure 2.14 presents the recovered delay vs. temperature result under *passive recovery* (0V) and *active recovery* (-0.3V) condition; in both cases, high temperature accelerates recovery. The proposed model accurately predicts this behavior, as also shown in the figure.



Fig. 2.13 Negative voltage-enabled active recovery after being stressed for 24 hours (Net delay increase is $\sim 3.24ns$). X-axis is the recovery time. (a) at 20°C, (b) at 110°C.



Fig. 2.14 High Temperature-accelerated recovery after being stressed for 24 hours (Net delay increase is $\sim 3.24ns$). X-axis is the recovery time. (a) under 0V (*passive* recovery), (b) under -0.3V (*active* recovery).

2.7.3 Model Validation

Figure 2.15 shows the measured raw data (delay change ΔT_d) vs. model predictions over time for four recovery cases and indicates that test results match the modeling results well, the analytical model provides an accurate estimation of the recovery under different conditions. Table 2.4 summarizes the recovery percentage (Recovered delay/Net delay increase) in four test cases. Both model prediction and measurement results are included. It is worth to mention that in the accelerated & active recovery case when both high temperature and negative voltage are applied to the system, about 72.4% of the degradation can be recovered within only 1/4 of the stress time, and this also indicates that high temperature and negative voltage can "assist" each other during recovery and lead to the maximum recovered portion (72.4% >> 16.7% + 28.7%). The recovered portion can be directly translated to the necessary design margin reduction. For example, with the accelerated & active recovery techniques, the design margin can be brought back to 27.6% of the original one within only a short period of recovery time. Figure 2.16 shows the measured frequency over the whole period of wearout and accelerated self-healing behaviors under high temperature (110°C), negative voltage (-0.3V) and stress vs. recovery ratio of 4. In summary, accelerated selfhealing techniques for BTI are effective compared to the passive recovery condition, but if we continue the accelerated self-healing period after 12 hours in the test case shown in Figure 2.16), BTI wearout can't actually be fully recovered even under the accelerated and active recovery conditions due to the existing of permanent components. In the following sections, we will explore this in depth.

Table 2.4 Summary of t	ne Accelerated	Self-healing	results	for 6	hours	of reco	very	(%:
recovered percentage)								

Test Case	Sleep Condition	Measurement Results	Model Prediction
Passive Recovery	20°C and 0V	0.66%	1%
Active Recovery	20°C and -0.3V	16.7%	14.4%
Accelerated Recovery	110°C and 0V	28.7%	29.2%
Accelerated Active Recovery	110°C and -0.3V	72.4%	72.7%



Fig. 2.15 Delay change (ΔT_d) over time during recovery: Model predictions vs. Measurement. Solid lines are model predictions, markers are measured data.



Fig. 2.16 One test case showing that accelerated self-healing techniques can recover about 72.4% of BTI wearout within only 1/4 of the stress time. Passive recovery data are not shown in the figure, but it is about less than 2% of recovery after 12 hours.

2.8 Reversible vs. Irreversible BTI Wearout

The accelerated self-healing techniques are able to rejuvenate the "aged" chip from BTI significantly. While based on the physical trapping and de-trapping mechanisms of BTI discussed in Section 2.2, and also validated by our experiments, there are still irreversible components (e.g. shown in Figure 2.3 as V_{th} net increase) that are hardly recovered within

a reasonable time and will accumulate cycle by cycle, thus hurting the performance and increasing the guardband. In this section, we look into the permanent component of BTI more in details and explore the frequency dependency of wearout and recovery. The same experimental setup and test methodology in Section 2.5 are also employed in this part of work.

2.8.1 Fast Traps vs. Slow Traps — A Physics Perspective

The BTI mechanisms [37, 135] suggest that the charge carriers need to gain sufficient kinetic energy to overcome a potential barrier necessary to break an interface state to be captured in the traps during trapping process. Here we define *fast traps* as those traps that have a high probability of trapping the charge carriers. These traps have a relatively lower trap energy barrier and are easier to be filled in a shorter time. On the contrary, *slow traps* are those have a higher trap energy barrier that is difficult for charge carriers to overcome. The principle of physics for de-trapping (recovery process) is that the trapped charge carriers (e.g. electrons for NMOS) have a certain probability to escape, with the probability being higher if their energy is higher and the trap energy barrier is lower, and vice-versa. Based on the statistical mechanics theory, the distribution of kinetic energies is proportional to the product of density of state and the Boltzmann distribution [62]. The 3-dimensional density of state is proportional to \sqrt{E} , therefore the energy distribution of electron is given by:

$$f_E(E) = A \times \left(\frac{1}{kT}\right)^{3/2} \times \sqrt{E} \times exp\left(-\frac{E}{kT}\right)$$
(2.18)

where *E* is the is the energy of the electrons, *k* is Boltzmann constant, *T* is temperature in *Kelvin* and *A* is a normalization factor. Since the de-trapping rate is proportional to the number of electrons at the energy of consideration, the de-trapping rate is proportional to $f_E(E)$. The energy distribution of electrons at room temperature (300*K*) is plotted in Figure



Fig. 2.17 Energy distribution of electrons at room temperature

2.17, which shows that majority of the electrons are at low energy in *meV* range, whereas the center energy of even the lowest energy of the trap is in order of several kT (~ 0.026*eV*) [88, 171]. This means that only a small fraction of electrons at the tail of the distribution could participate in the de-trapping process.

Shown in Figure 2.18 is the illustration of the trapping and detrapping process for two types of the traps. Since fast traps have lower trap energy barrier, so it is easier for charge carriers to escape from them, and this leads to *fast recovery*, or reversible part of wearout. For the slow traps, the charge carriers need to overcome a higher trap energy barrier. As a result, it is very slow or even impossible for them to escape within a reasonable time, so these traps will cause the *irreversible wearout*. To give a first order estimation, if the trap energy is 100meV higher, the probability distribution goes down by a factor of $exp(100meV/20emV) \approx 50$. This indicates that the time it takes to de-trap goes up by about a factor of 50 for every 100meV increase in trap energy.



Fig. 2.18 Illustration of Fast traps vs. Slow traps.

Temperature and voltage (electrical field) play an important role of determining energy of electrons. Figure 2.19 shows that by increasing the temperature, the energy distribution shifts to the right, so the probability of de-trapping increases. This indicates that the boundary between the reversible and irreversible wearout is not fixed and can actually even be shifted.

2.8.2 Irreversible Wearout during Accelerated Self-healing

As discussed in Section 2.7.3, although the part of the irreversible wearout for passive recovery could be recovered by accelerated active recovery, and this has been demonstrated by our experiment, the measured result is shown in Figure 2.20. The frequency under passive recovery condition (27°C and $V_{gs} = 0V$) is normalized to the accelerated self-healing condition. It is clear that in both cases, recovery saturates to some values below fresh state frequency, and the irreversible part of BTI wearout after passive recovery is much larger than the accelerated & active recovery case.



Fig. 2.19 Energy distribution of electrons at room temperature.

The "unchecked" irreversible components will keep accumulating throughout the system lifetime. Figure 2.21 is a test case for several cycles of stress and recovery. In each cycle, a 6-hour accelerated stress is followed by a 6-hour accelerated recovery (6 hours vs. 6 hours). *IRn* refers to the amount of irreversible wearout for the *nth* cycle. It shows in the figure that the recovery under accelerated conditions saturates in each cycle, and the irreversible wearout (IR) increases for the first few cycles and settles afterwards. A possible explanation for this behavior is that in the later cycles, some of the irreversible wearout from previous cycles starts to recover, and the accelerated recovery and stress can fully compensate each other. But it will not be fully recovered to the fresh state even with accelerated active recovery techniques applied during each cycle.

2.8.3 Sequentiality of Reversible and Irreversible Wearout

As reversible wearout and irreversible wearout are mostly determined by *fast* traps and *slow* traps respectively. So there will be sequences when wearout happens due to the different



Fig. 2.20 Irreversible part under two recovery conditions – Passive recover vs. Accelerated & active recovery (accelerated self-healing). BTI doesn't recover completely even under the accelerated self-healing conditions.



Fig. 2.21 Accumulation of the irreversible BTI wearout cycle by cycle. IRn refers to the accumulated irreversible component after nth cycle.
trapping rate of the two. To further investigate this, a group of stress tests is conducted. Shown in Figure 2.22 is the frequency degradation under two accelerated stress conditions with different stress voltages. It illustrates that both test cases follow similar wearout patterns. Firstly, reversible wearout kicks in, then the effect of reversible wearout levels off and irreversible wearout takes over – in time domain this is roughly seen as a steep slope followed by shallow slope during wearout.

Figure 2.23 shows a test case when a 6-hour accelerated & active recovery (110° C, -0.3V) durations follows a 6-hour stress. The accelerated recovery speeds up the recovery process and even recovers some parts that would otherwise be considered irreversible. In the time domain this can be roughly seen as a steep slope followed by a *saturation* (zero slope) once all the reversible part and part of irreversible part were recovered. Also, it is worth to mention that the *Recovered Wearout* is larger than the *Reversible Wearout* as shown in the figure, and this further demonstrates that the *accelerated & active recovery* techniques are able to recover some of the irreversible parts. But solutions that are able to fully fix or avoid this component are still necessary and highly preferred, they are presented in the next section.

2.9 Frequency Dependency of BTI Wearout and Recovery

2.9.1 Sleep with Accelerated Rejuvenation when Getting Tired

The whole process of wearout and accelerated active recovery can be compared to the biological world. Humans for example are active during daytime, with the body conducting activities and experiencing fatigue. During night time sleep, the body goes through several active processes that are essential for the recovery of its full capabilities for the next day. If some organs experience heavy fatigue without in-time rest, part of the fatigue will be translated into some potential harms to the body, and will be hardly recovered. This is well



Fig. 2.22 BTI-induced frequency degradation under two accelerated stress conditions. In both cases, reversible wearout kicks in firstly, then it levels off and irreversible wearout takes over.



Fig. 2.23 Sequentiality of reversible and irreversible wearout.

known for athletes that need scheduled recovery periods after extensive workouts, with their athletic performance actually getting better after the rest periods. These biological fatigue and recovery schedules are not unlike those illustrated in Figure 2.21. We thus borrow these ideas and see how they apply to electronic systems. The key idea is to stop the stress before irreversible effects get a chance to accumulate. The ideal strategy is thus to keep the circuit active only during the reversible phase of wearout until the irreversible wearout kicks in; thus the irreversible wearout becomes almost unobservable even in accelerated stress cases. To validate this idea, a set of tests with different "circadian rhythms" is conducted, and is summarized in Table 2.5. All tests start from the fresh state. The total test time is 3 days for all test cases.

Table 2.5 Summary of periodic accelerated rejuvenation test cases

Test Case	Chip No.	Cycle Stress Time	Cycle Accelerated & Active Recovery Time	# of cycles
6 hrs vs. 6 hrs	1	6 hours	6 hours	6
4 hrs vs. 4 hrs	2	4 hours	4 hours	9
2 hrs vs. 2 hrs	3	2 hours	2 hours	18
1 hr vs. 1 hr	4	1 hour	1 hour	32

Note: All accelerated & active recovery are under -0.3V and 110° C.

To make a fair comparison of recovery percentage among different chips and also under different operating conditions, normalization is necessary. From Equation 2.4 and 2.5, we can estimate the change in gate delay Δt_d when V_{th} is subject to change as:

$$\Delta t_d \sim \Gamma \cdot \Delta V_{th} \cdot t_{g0} \tag{2.19}$$

where t_{g0} is the time zero delay of the gate with no V_{th} shift, Γ is a constant. Assume that Chip 1 has an initial frequency (fresh status) of $f_{1TA}(0)$ at temperature T_A , and chip 2 has an initial frequency of $f_{2TB}(0)$ at temperature T_B . If two chips undergo the same threshold voltage change $\Delta V_{th}(t)$ after being stressed or recovered for time length of t, and since the temperature induced threshold voltage shift doesn't change with time, so the resulted frequency of Chip 1 at temperature T_A becomes:

$$f_{1TA}(t) \sim \frac{1}{t_{1TA}(0) + \Gamma \cdot \Delta V_{th}(t) \cdot t_{1TA}(0)}$$
(2.20)

Similarly, the frequency of Chip 2 at temperature T_B is:

$$f_{2TB}(t) \sim \frac{1}{t_{2TB}(0) + \Gamma \cdot \Delta V_{th}(t) \cdot t_{2TB}(0)}$$
 (2.21)

The process of normalizing Chip 2 data to Chip 1 data is done through

$$\frac{f_{2TB}(t)}{f_{2TB}(0)} \cdot f_{1TA}(0) = \frac{t_{2TB}(0)}{t_{2TB}(0) + \Gamma \cdot \Delta V_{th}(t) \cdot t_{2TB}(0)} \cdot \frac{1}{t_{1TA}(0)}$$

$$= \frac{1}{t_{1TA}(0) + \Gamma \cdot \Delta V_{th}(t) \cdot t_{1TA}(0)} = f_{1TA}(t)$$
(2.22)

The above derivation process proves that although the same amount of BTI-induced threshold voltage $\Delta V_{th}(t)$ could lead to different frequency degradation for different chips under different temperatures, the normalization process could reflect the actual threshold voltage shift across multiple chips properly. Thus, in the following measurement results section, we are able to normalize the frequency of all 4 test cases shown in Table 2.5 to one of them to make fair comparisons.

2.9.2 Measurement Results

Shown in Figure 2.24 are the measurement results. For all test cases except the 1 hr vs. 1 hr case, the accelerated active recovery has a period of saturation which indicates the irreversible parts of the wearout, and the irreversible parts accumulate in the first several cycles and settle down in the following cycles. For the 1 hr vs. 1 hr case, alternating phases of stress and accelerated recovery can completely compensate for each other, and after each accelerated active recovery phase the chip can indeed start fresh. *The irreversible part*



Fig. 2.24 Measurement results for different "circadian rhythms".

of wearout is totally avoided explicitly. Figure 2.25 presents the accumulated irreversible wearout for the first 6 cycles under above four test conditions. It shows that the earlier the accelerated rejuvenation techniques are applied, the slower the irreversible wearout accumulate, which results in less permanent component. There is an optimal balance of



Fig. 2.25 Irreversible component accumulated during the first 6 cycles for four different scheduled stress and accelerated and active recovery shifts (circadian rhythms).

stress and accelerate recovery (e.g. 1hr vs. 1hr in this accelerated case) which leads to almost no irreversible wearout.

Assume that the amount of frequency degradation under normal condition is the same as the accelerated case for 1 hour, the equivalent duration is about *31 hours* under nominal voltage and room temperature based on the model proposed in Section 2.4.2. As shown in Figure 2.26, the identical optimal condition could be that the chip is active under normal operating conditions for (at most) *31 hours*, the accumulated BTI wearout could be fully mitigated by a following *1 hour* (or longer) accelerated active recovery duration.



Fig. 2.26 An identical regular-operation use case (31 hrs vs. 1 hr) to the 1 hr vs. 1 hr accelerated stress case for FULL recovery.

2.9.3 Reduction of Necessary Design Margin

The explored unique behaviors of BTI wearout and (accelerated & active) recovery provide the big potentials of reducing the necessary design margins significantly during the early design phase. As discussed in Chapter 1, for the worst case design solution, to meet the timing requirement throughout the whole lifetime, guardbands need to be added (e.g. by oversizing transistors). Without the proposed periodic accelerated and active rejuvenation solutions, the margin needs to cover both reversible and irreversible wearout, and the irreversible part has to cover long time periods (typically years). Since the proposed strategy starts recovery before the irreversible wearout kicks in, the design margin only need to cover reversible wearout. Assume that the irreversible wearout at room temperature is the same as the one at the accelerated stress case (at 110°C), and we define AC stress as the case when transistors switch between ON and OFF with a 50% duty cycle, which gives the balanced stress and recovery during operation. Based on the model in Section 2.4.2, the estimated design margin of 5 year and 10 year lifetime spans under DC and AC stress at room temperature is shown in Fig. 2.27. The proposed solution (1hr. vs. 1hr case) gives a



Fig. 2.27 Necessary design margin estimation under different stress conditions.

design margin reduction of at least $60 \times$ for all cases. In the AC stress case for a 10-year lifetime constraint, the design margin reduction is more than $100 \times$. It also shows that as the lifetime constraint increases, the guardbands need to be relaxed (2×) correspondingly, while with the proposed strategy, the design margin stays almost the same (1×).

2.9.4 Reduction of Tracking Power

The alternative solutions for dealing with wearout are adaptive techniques at the circuit level [189] or dynamic reliability management techniques at the architecture level [199], where wearout sensors are deployed to track during the whole period of the lifetime. This means more tracking power. With the proposed strategy, the time for recovery is known ahead, wearout sensors only need to track during a short time (e.g. 31 hours shown in Section 2.9.1) for the reversible part of wearout. Numerically the difference for the tracking power

is between O(ln(31hours)) and O(ln(10years)) for a 10-year lifetime constraint, or about $ln(2826) \sim 8 \times$ reduction.

2.9.5 Average Performance Improvement

With the periodic accelerated and active rejuvenation, the circuit is guaranteed to run faster than the case when no recovery is applied. We define that the average performance refers to the average of all performance values during operation time (wearout period). Figure 2.28 shows the average performance improvement (IMP) calculated from the measurement results for 1 day and 2 days with the same chip. Under the 1hr vs. 1hr case, when irreversible wearout is almost completely avoided, it gives the best average performance, which is close to the fresh status. As operation time increases (e.g. 1 day to 2 days), the average performance will keep almost the same for 1hr vs. 1hr case, while for other test cases, especially the case when no recovery strategy is applied, the average performance decreases dramatically, and this leads to that the average performance improvement achieved by the proposed strategy will increase with time $(1.6 \times \text{ from 1 day to 2 days)}$.

Figure 2.29 presents the predicted average performance improvement over the no-recovery solution under nominal operation conditions (room temperature, nominal voltage) predicted by the model. As the lifetime constraint increases from 5 years to 10 years, the average performance for the proposed rejuvenation solution will keep almost the same, compared to the no-recovery case when the average performance scales down dramatically. In other words, the average performance improvement enabled by the proposed solution will increase significantly as the lifetime requirement increases. To give an example, for a 10-year lifetime span, the improvement can be as large as $\sim 9\%$.



Fig. 2.28 Average performance improvement (IMP) for 1 day and 2 days from measurement.



Fig. 2.29 Average performance under different stress conditions.



Fig. 2.30 Frequency dependency of irreversible component of BTI wearout.

2.9.6 Frequency Dependency Behaviors of BTI Wearout

It has been shown in many literature [57, 42] that BTI-induced V_{th} shifts due to AC stress are independent of frequency in the range of few Hz to GHz, and most of the BTI models are also developed based on this theory, while our experimental results demonstrate that independence of BTI on frequency doesn't hold for the whole frequency spectrum. Figure 2.30 shows the measured permanent component vs. frequency under accelerated stress and recovery conditions. As in the higher frequency range (close to 0 on the X-axis), the permanent component is almost zero and doesn't increase with frequency. When the frequency reaches the lower range (moving to the right on the X-axis), the accumulated permanent component increases inversely with frequency. There is a turning point (range) where "defines" the boundary between frequency dependency and independence. To the best of our knowledge, this thesis is the first to study experimentally the BTI effect under very slow AC stress with the frequency range of 1/hours and also under accelerated stress and

recovery conditions, the explored frequency dependent behaviors can potentially lead to that the permanent component can be minimal through proactively scheduled recovery intervals. Implementation details will be further discussed in Chapter 5.

2.10 Conclusions

As BTI becomes one of the dominant reliability challenges for present and future digital circuits, most of the previous BTI mitigation techniques focus on reducing BTI-induced degradation during operation (under stress) or utilize the passive recovery behavior of BTI, however BTI is not fundamentally repaired in those cases, and this could cause permanent failures. In this chapter, we first presented a series of accelerated self-healing techniques which are able to "reverse" the direction of wearout, thus accelerating and activating the BTI recovery. Based on the actual hardware measurement results with 40nm FPGAs, these accelerated self-healing techniques can lead to significant recovery rate improvement (e.g. we demonstrated a case where 72.4% of the wearout is recovered within only 1/4 of the stress time for BTI). While even in the accelerated self-healing recovery conditions, there are still part of BTI wearout that are irreversible. We further explored the boundary between reversible and irreversible wearout physically and experimentally. The main findings are: First, we show that the boundary between the reversible and irreversible parts of transistor wearout is not fixed, with the irreversible part becoming at least partially reversible under the right conditions of accelerated active recovery and stress/recovery scheduling. Second, we show that there are certain stress/recovery schedules that can (almost) completely eliminate irreversible wearout, thus allowing significant reductions in necessary design margins $(>60\times)$ and improvement in average performance ($\sim 9\%$ with a 10-year lifetime). The discovered BTI frequency dependency is able to help the community understand BTI behaviors more deeply, these unique BTI recovery properties introduce a new knob for designing reliable systems. Figure 2.31 highlights the key contributions of this chapter. The potential implementations and



Fig. 2.31 Chapter 2 highlights.

tradeoff analysis by taking use of the accelerated self-healing techniques will be discussed in the following chapters.

The work presented in this chapter has been published in [J2], [C12] and [C14]. It was also presented in [P5], [P6], [P11] and [T4].

2.11 Acknowledgments

I would like to thank NSF (Grant No. CCF-1255907) and SRC Global Research Collaboration (GRC) Program (Task ID. 2410.001) for sponsoring this work. I also thank Prof. Wayne Burleson (University of Massachusetts, Amherst) and Alec Roelke (University of Virginia) for discussions and feedback.

Chapter 3

Accelerating and Activating Recovery for EM Wearout

3.1 Overview

The down-scaling of CMOS technologies into the nano-regime not only elevates the transistor wearout issues such as BTI, it also worsens the interconnect (on-chip metal wire) wearout effect Electromigration (EM). Especially, EM has been a significant concern in power delivery networks (PDN), which largely experience unidirectional current flow [173, 163]. EM-induced failure is projected to get even worse due to the increasing current densities from shrinking interconnect geometries in the sub-10nm regime [66]. EM occurs due to the gradual displacement of metal atoms in a semiconductor. When the current density is high enough, it can cause the drift of metal ions in the direction of the electron flow. As BTI degrades chip performance by slowing down device switching speeds. EM can increase wire resistance, which can cause voltage drop resulting in device slowdown; it can also cause permanent failures in circuits due to shorts or opens if the stress accumulates to certain amount. Conventional ways of addressing EM effects are by design rules (e.g. metal width requirement) during the physical design phase and signoff phase [126]. While Conservative



Fig. 3.1 Illustration of EM Recovery: Stress relaxation occurs when current is switched off, and this is similar to the passive recovery condition for BTI.

oversizing the metal can significantly sacrifice area, power and routing costs. In addition, the dynamic workloads and changing operating conditions can still lead to large variations in current density, which can cause major EM threats [109].

Similar to BTI, one important transient behavior of EM wearout is the recovery property, which refers to the stress relaxation in the metal line when the current is switched off as shown in Figure 3.1. This can refer to the passive recovery condition as in BTI recovery case. Under this condition, the effect of the electron wind induced-stress can be relieved to only certain levels as demonstrated with experiments in [118, 128], but it can not be fully released due to the existing of permanent component. In this chapter, we demonstrate several solutions that can activate and accelerate the recovery of EM experimentally with a group of on-chip metal lines, these solutions are inspired by the notion of "reversing the direction of wearout". We explore that EM wearout also has frequency dependent behaviors where the amount of permanent component depends on the periods of stress and recovery (with the same duty cycle), not unlike BTI. By inserting the accelerated and active EM



Fig. 3.2 Electromigration Mechanism: EM is the result of the dominant force F_{wind} , that is, the momentum transfer from the electrons which move in the applied electric field.

recovery periods periodically, the overall mean-time-to-failure (MTTF) of metal line can be significantly extended. Since during EM active recovery, the current still flows, thus the system doesn't need to be switched off, this can lead to minimal performance overhead if implemented with the necessary assist infrastructures. This chapter presents the experiment setup, theory and measurement results for the proposed active and accelerated EM recovery techniques.

3.2 EM Wearout and Recovery Mechanisms

Compared to BTI wearout, EM mechanism is less debatable. It is widely accepted that EM is the result of momentum transfer from the electrons, which move in the applied electric field, to the ions which make up the lattice of the interconnect material. Figure 3.2 illustrates the process. Current flow through the metal line produces two forces, the first one is electrostatic force F_{field} caused by the electric field strength in the metallic interconnect, and the force is usually small and can be ignored. A second force F_{wind} is caused by the momentum transfer between electrons and metal ions in the crystal lattice. This force acts in the direction of current flow and is the main EM source [126]. Over time, this can lead to resistance increase and open circuit.

The healing effects (EM Recovery) refer to those caused by the atomic flow in the direction opposite to the electron wind force F_{wind} , the back-flow, during or after EM. This



Fig. 3.3 Electromigration Stress and Passive Recovery: EM mainly affects the power delivery network (PDN). EM healing occurs when current is removed, but the recovery is partial and slow.

back-flow of mass begins to take place once a redistribution of mass has begun to form. It tends to reduce the failure rate during EM and partially heals the damage after current is removed. The cause of this back-flow of mass is the inhomogeneities, such as temperature and/or concentration gradients, resulting from EM damage [49]. Due to this effect, the signal metal lines suffer less EM effects because most interconnects are stressed under bidirection current which correspond to the charging and discharging processes. While for power delivery network (especially global PDN) shown in Figure 3.3, EM doesn't exhibit the luxury of being recovered because of the unidirectional current flow, and this can cause huge power net IR-drop and permanent failures potentially. Thus, in this work, we mainly look into the EM issues in power delivery network.

EM is characterized by mean time to failure (MTTF) conventionally. The MTTF of a single metal interconnect caused by EM is given by well-known Black's Equation [25],

which is:

$$MTTF = \frac{A}{J^n} exp(\frac{E_a}{kT})$$
(3.1)

where A is a constant depending on the cross-sectional area of the interconnect, J is the current density, n is a scaling parameter, which usually equals to 2 for voild-nucleation-limited failure and 1 for void-growth-limited failure [75]. k is the Boltzmann constant, T is the temperature in *Kelvin*. E_a is the activation energy. Equation 3.1 shows that current density J and the temperature T are deciding factors that affect MTTF due to EM.

Recent EM modeling framework [84, 85] has shown that EM wearout experiences several phases. The first phase is called void nucleation phase during which the stress accumulates until it reaches the critical value $\sigma = \sigma_{critical}$, the resistance during this phase is almost unchanged; Following the void nucleation phase, the generated voids start growing and lead to an increased resistance over time. As a result, the PDN becomes a time-varying network and the voltage drop changes over time.

3.3 **Prior Work on EM Recovery**

The recovery effect of EM under AC stress was firstly studied in [209]; the experimental results showed that the EM lifetime increases with frequency. This effect was further analyzed in [2], which demonstrated that the healing can increase the lifetime by several orders of magnitude depending on the metal used. In [222], authors looked into EM recovery in TSVs of 3DIC. While [118] suggested that EM is not fully recovered even during an opposite polarity pulse current with 50% of duty cycle, this means that EM stress and passive recovery are not symmetric, and there is an irreversible component for EM as well. Modeling work [85, 204] suggested from a physics perspective that high temperature can speed up atom diffusion towards the cathode end and lead to faster and more recovery, but these work are still simulation based, there are no experimental results are presented. [155] was the first to

conduct EM recovery experiments under different temperatures, but the goal of that work was to understand the transient resistance change during a temperature sweep up to 400*K*, and a first-order model was proposed to capture the behaviors. In [174], authors demonstrated the argument that at low frequency, when current flow is interrupted, the stress gradient is sufficient to effectively counter the effect of EM and allow stress relaxation and consequently longer lifetimes. But it was still under passive recovery where the current is turned off in pulsed DC (PDC) operations. In our work, we explore both temperature and current direction impacts on EM recovery, we also demonstrate the frequency dependency behaviors for EM wearout. Through multiple "deep healing" methods, EM-induced MTTF can be significantly extended.

3.4 "Reversing" the Direction of EM Wearout

In this section, we propose that EM recovery can be made active and accelerated. During passive recovery period, several boosting techniques can be applied, and they are shown in Figure 3.4, condition No. 2 shows the active recovery case where the direction of current is reversed to activate and assist the electron back flow. In last chapter, we discussed that high temperature can increase the kinetic energy for the charge carriers. More importantly, increased temperature can lead to recrystallization accompanying an increase in grain size and defect decay, such as the annihilation of vacancies at film surfaces [155]. Overall, the EM-induced damage and stress can be healed and relaxed by high temperature, and this case is shown as No. 3. No. 4 illustrates a combined condition of No. 2 and 3. To validate these recovery conditions, we also conduct actual hardware testing on a set of on-chip metal lines. Details of setup is presented in the next section. Since EM permanent component hasn't been studied and well understood, we also study the EM frequency dependent behaviors with the same setup.



Fig. 3.4 Electromigration "reversing" techniques: No. 1 refers to the passive recovery, No. 2, 3 and 4 are proposed active and accelerated recovery solutions.

3.5 Test Setup

3.5.1 Test Structure

Since EM mainly happens to on-chip metal wires, and there is no commercial test infrastructures which contains only on-chip metal lines. We fabricated a test chip with 180*nm* bulk CMOS technology using dual-damascene Cu interconnect, the test structures are a series of "long" and "narrow" on-chip metal lines. Figure 3.5 shows the die photo along with the dimension of the metal wire. The metal wire is fabricated with the highest metal layer (M6) of the technology in dual-damascene process since this metal layer is more likely being used in the global power delivery network. The resistance change ΔR over time is measured during stress and recovery phases. Based on Equation 3.1, EM depends on the



Technology	180nm
Material	Copper
Thickness	0.8um
Length	2.673mm
Width	1.57um
Resistance (@rt)	35.76 Ω

Fig. 3.5 Die photo with the test structure for EM recovery: On-chip "long" and "narrow" metal lines and their dimensions; "rt" means room temperature ($\sim 27^{\circ}$ C).

current density flowing through the metal line, which is inversely proportional to the cross section area, since the metal thickness depends on the process itself, so we can only control the width, which is picked being "narrow" enough $(1.57\mu m$ for the 180*nm* technology) while not violating the design rules. Although Equation 3.1 doesn't suggest any EM dependency on metal length, shorter metal experience less or no EM due to immortality condition which is also known as Blech limit [26]. According to the theory, the immortal metal segments can be filtered based on the following condition:

$$(j \times l) \leq (j \times l)_{critical} = \frac{\Omega \sigma_{critical}}{eZ\rho}$$
 (3.2)

where *l* is the metal length, Ω is the atomic volume, *e* is the electron charge, *eZ* is the effective charge of the migrating atoms, ρ is the wire electrical resistivity, $\sigma_{critical}$ is the critical stress needed for the failure precursor nucleation (void/hillock). Equation 3.2 means that the length of the metal line needs to be picked long enough to capture the EM prominence. As shown in Figure 3.5, our metal line test structure is across the whole die area, and the length is 2.633*mm*. Probe pads are used for wire bonding to connect with the external measurement board, they can also be probed directly.

3.5.2 Measurement Setup



Fig. 3.6 Electromigration stress and recovery measurement setup.

In our measurement setup, we first bond the wire-under-test to a regular Dual in-line (DIP) package. Compared to the on-chip metal wire, the bonding wires employ a much larger cross-sectional area (> 10×), and are believed to be much less impacted by EM. So the overall EM effects are dominated by on-chip metal wire. Figure 3.6 shows the whole measurement setup, where the bonded chip is connected with a constant current supply using high temperature wires (red and black wires shown in the figure) on a breadboard. In this way, only wire-under-test is exposed to the high temperature environment. A Voltmeter (Analog Discovery from Digilent [12]) is connected in parallel to the test structure for measuring the voltage drop ($\Delta V(t)$) on the wire-under-test. The voltmeter is USB-powered and supports visualization and data recording. We sample the voltage every minute and saved it in a file

for further processing. The resistance change over time due to EM can be calculated as $\Delta R(t) = \Delta V(t)/I$ based on the Ohm's law, where *I* is the constant current fed to the wire. The current value *I* and stress temperature *T* are decided based on both Equation 3.1 and exercise experiments so that MTTF is estimated to be in the range of few days. Before stress phase, we first wait for a reasonable amount of time to ensure that a steady-state temperature was reached by the thermal chamber which allows temperature fluctuation of $\pm 0.3^{\circ}$ C only. The current supply is able to provide bidirectional DC current, and this allows a short switching time (less than 2*s*) between stress and recovery phase.

3.5.3 Test Cases

Phase	Case Index	Wire No.	Τ (° C)	$\mathbf{J} (MA/cm^2)$	Time (hours)	Comments
	EMST1	1	27	7.96	12	Baseline for Stress
Stross	EMST2	2	230	7.96	12	-
(Activo)	EMST3	3	230	7.96	20	Metal broke
(Active)	EMST4	4	230	7.96	12	-
	EMST5	5	230	7.96	6.7	-
Docovory	EMPR1	2	27	0	20	Baseline for Recovery
(Sloop)	EMASH1	4	230	-7.96	10	-
(Sieep)	EMASH2	5	230	-7.96	10	-

Table 3.1 Test Cases for EM Stress and Recovery

All accelerated stress tests are conducted on "fresh" chips with the same test structures. Table 3.1 summarizes all the test cases, where "EMST" stands for accelerated stress conditions, "EMPR" refers to passive recovery and "EMASH" represents accelerated and active recovery conditions. The measurement halts when EM breakdown occurred for example in case EMST3, which provides a reference on how long it takes to reach the metal breakdown point (TTF). EMST1 case is used as a baseline stress condition for comparisons. Similarly, EMPR1 case is used as a baseline recovery condition for comparing against the accelerated and active recovery conditions (EMASH1 and EMASH2). EMST4 and EMST5 are stressed under the same accelerated conditions except that EMST4 is stressed with a shorter

period. These two stress periods are followed by two accelerated and active recovery periods (EMASH1 and EMASH2) respectively.

3.6 Experimental Results for EM Active and Accelerated Recovery

This section presents the experimental results from test cases summarized in Table 3.1. Shown in Figure 3.7 is the measured EM-induced resistance change under accelerated stress and recovery conditions with relatively high constant current density $(\pm 7.96MA/cm^2)$ and elevated temperature (230°C). During the accelerated stress phase, the results indicate that the EM evolution consists of two distinct phases that were described in Section 3.2 – the void nucleation phase and the void growth phase. During the nucleation phase, the EM-induced stress increases until it hits a critical value, when voids are generated; before this point, the resistance has almost no change. Following the void nucleation phase, these generated voids start growing and lead to an increased resistance over time. Our experimental results agree with measured data in [201, 209], and are also consistent with what is predicted by recently proposed physics-based EM models [86, 204].

During the active and accelerated recovery phase, a reverse current (with the same absolute value as in the stress phase) and elevated temperature are applied; Figure 3.7 shows that the activated recovery is much faster than that under passive recovery, and more than 75% of EM wearout can be recovered within only 1/5 of the stress time. Figure 3.7 also shows the results under passive recovery condition, where recovery saturates to a very high-resistance value after a short period of true recovery, and this saturation continues even with extended recovery periods.

However, our results also suggest that there is still a lingering permanent component in accelerated and active recovery case, which has similar behavior to what we saw for initial



Fig. 3.7 Measurement results for EM degradation and recovery under passive recovery (Figure 3.4 No. 1 or Test case EMPR1 in Table 3.1) and proposed recovery conditions (Figure 3.4 No. 4, at 230°C and $\pm 7.96MA/cm^2$) during the void growth phase: there is still a permanent component even under accelerated and active recovery.

BTI wearout measurements. This inspires us to explore the frequency dependent behavior that is not unlike the method we used in BTI recovery case. Applying in-time recovery for EM enables the hope of reducing, avoiding or even eliminating the permanent component of EM; Figure 3.8 demonstrates exactly this. The results show that by scheduling the recovery phase in the early phase of void growth, EM can also be fully recovered. But the potential issue of scheduling recovery during void growth is that during recovery, there is still (reverse)



Fig. 3.8 Measurement results for EM accelerated and active recovery during the early period of the void growth phase (at 230°C and $\pm 7.96MA/cm^2$): full recovery.

current flowing through the metal, and this could lead to potential EM, but in the opposite direction (shown in the figure), and thus add uncertainties in terms of ultimate effects. A more "economic" way is to schedule the recovery periodically before voids nucleation happens; the measurement results of this strategy are shown in Figure 3.9, where multiple short accelerated and active recovery intervals (30*min*) are inserted in the early phase of EM stress evolution, and this results in a delay of void nucleation for a significant amount of time (almost $3 \times$ slower compared to Figure 3.7). By employing such circadian rhythm-like scheduling recovery strategy, the overall time-to-failure (TTF) can be potentially significantly extended. For example, Figure 3.9 shows TTF is increased by more than $2 \times$ with 5 short recovery periods (2.5hours in total) scheduled in the early lifetime.



Fig. 3.9 Measurement results for scheduled periodic recovery intervals during void nucleation phase: It takes much longer for voids to nucleate, and the overall TTF is extended significantly.

Based on extensive accelerated stress and recovery tests, we conclude that EM recovery (back-flow) can be further activated and accelerated significantly, the "Push-Pull" stress/active recovery compensation where in-time scheduled periodic recovery intervals are able to fully eliminate the permanent EM component. The insertion of short intervals of EM recovery before void nucleation is able to extend the lifetime of the metal wires significantly. While BTI active recovery needs to be in an OFF period, and EM active recovery happens during ON period when there is reverse current flowing; this opens new opportunities of scheduling both recovery over the whole lifetime span with the proper circuit solutions, which will be discussed in details in the following chapters.

3.7 EM Signoff Considering Accelerated and Active Recovery

EM issues are usually addressed by back-end-of-line (BEOL) design rules or/and by adding design margins. In the first method, current rule limits in metal wires are set by foundry to ensure reliable operations over a prescribed time period without significant EM damage. Designers need to follow these rules when doing floorplanning and physical design. For example, the number of power straps, number of VIAs and PADs will be decided partially by these EM design rules. EM signoff tools (such as Cadence Voltus [164]) can also be employed to analyze and optimize the full-chip EM and IR drop. In the second method, design teams are forced to use larger margins to guardband against EM-induced delay degradation, that is, design and sign off at either a shorter lifetime or slower speed [40]. This margin could also be voltage margin which is added when the resistance of the PDN mesh increases due to EM to keep the drive current the same so as to achieve the same performance. In summary, both methods require extensive design-time estimations and lead to some overhead for some metrics such as performance.

With the explored accelerated and active recovery behaviors in this chapter, we show that EM can be fully recovered under certain conditions, the EM signoff requirements can be relaxed in several aspects. This section discusses three different scenarios that can take advantage of the EM recovery properties.

3.7.1 Relax the EM Design Rules

Assume that the EM recovery techniques have been implemented on-chip, and EM recovery can be accelerated and activated, full recovery can be achieved ¹. For ease of

¹Details of the implementing EM recovery circuitry on chip will be discussed in Chapter 4. Overheads and tradeoffs of adding these circuitry will be studied in Chapter 5. For this part, we assume this has been done and validated, and we focus on studying the potential benefits of EM accelerated self-healing techniques.



Fig. 3.10 Illustration of EM Current Rule Relaxation due to Recovery: Current density requirement can be potentially relaxed, x is a number larger than 1.

comparisons, we refer the regular case without any proactive recovery as Case EM_Reg, and the case with recovery techniques (accelerated self-healing) as Case EM ASH. If the lifetime target keeps the same and minimum metal width has been used, EM current design rule limits can be relaxed in Case EM ASH due to that in-time proactive recovery can always bring the "aged" metal line back to almost the original state, thus this will result in less required power straps and less routing congestions. This is illustrated in Figure 3.10. To quantify how much this current requirement can be relaxed (calculate the value of x), we refer an industry-standard 28nm FDSOI² technology as an example, since it is a relatively advanced node that has been deployed in many products, so we believe this analysis can be representative enough to cover the real design scenarios for modern chips. Listed in Table 3.2 are the maximum DC current allowed at junction temperature of 125°C for the top two metal layers (used for power routing) in this technology. The EM current rules assure reliable operations of 10 years. Since the current limit is given at 125° C, I_{dc} can be corrected by a temperature derating factor γ_{der} that is given in Table 3.3 for several temperatures. There is also a width derating factor that determines the maximum current, which can be relaxed by wider metal. Since in our case we assume the minimum width is used, so this width derating

²FDSOI refers to Fully Depleted Silicon On Insulator technology.

Metal Level	I_{dc} (mA)	Minimum Width (um)	Thickness (um)
M10 (Highest Level)	0.408	0.4	0.88
M9	0.408	0.4	0.88

Table 3.2 EM Line Current Limits @125°C for 10 years of operations

Table 3.3 Temperature Derating Factor $\gamma_{der}(T)$ and EM Line Current Limit

$T(^{\circ}C)$	$\gamma_{der}(T)$	$I_{dc}(T) EM(mA)$
125	1	0.408
110	2.792	1.139
50	40.62	16.573
27	68.312	27.871

Table 3.4 Summa	ry of Parameters	s for Black	s's Equation

Parameter	Value	Unit
k	$1.38 imes 10^{-23}$	J/K
E_a	0.9	eV
n	1	-
A	1.35×10^5	-

factor is 1. Overall, the maximum DC current required by EM rules at temperature T is given by:

$$I_{dc}(T)|EM = I_{dc}(125^{\circ}C) \times \gamma_{der}(T)$$
(3.3)

where $I_{dc}(125^{\circ}C)$ is listed in Table 3.2. The current limit at various temperature can be calculated based on Equation 3.3 and are listed in Table 3.3.

Based on Equation 3.1 and our accelerated testing results, we can get a first-order estimation of the stress and recovery under normal operation conditions (normal current and temperature). Assume in the normal case, the junction temperature is 50°C, which is close to a normal CPU die temperature after cooling. The current under this temperature is given in Table 3.3. The parameters used in Black's Equation are from the technology design rule manual, and they are listed in Table 3.4. We also assume that the same temperature and current values are applied during both stress and recovery periods. An equivalent stress and recovery conditions under normal operating condition to Figure 3.8 is that the voids

Description	Condition	Accumulation Time	Nucleation Time	Full Recovery Time	Comment
Accelerated	$J = \pm 7.96 MA/cm^2$	360 min	81 min	60 min	From
Condition ³	$T = 230^{\circ}\mathrm{C}$	500 min	01 <i>min</i>	0 <i>9 min</i>	Measurement
Normal	$J = \pm 4.71 MA/cm^2$	2.05 years	168 days	140 days	From
Condition	$T = 50^{\circ}\mathrm{C}$	2.05 years	108 auys	140 <i>auys</i>	Calculation
Relaxed	$J = \pm 9.42 MA/cm^2$	1 year	84 days	70 days	From
Condition ⁴	$T = 50^{\circ}\mathrm{C}$	i yeur	0+ uuys	10 auys	Calculation

Table	3.5	Estimated	EM	stress	and	recovery	time	under	normal	operating	conditions	for
28nm	FD:	SOI techno	logy									

accumulate for about 2.05 years (t_{nuc}) until nucleation, during this period the resistance keeps almost the same. Then the nucleation and growth periods start and continue for 168 days, the direction of the current is then reversed so that the active recovery starts, it will take about 140 days to fully recover back to the "fresh" state. This pattern can be replicated and continued. With this, the resistance will not reach the rapid void growth phase which can potentially lead to end of life. This opens the opportunities of relaxing the EM design rule. An example is illustrated in Figure 3.11. It shows two cases where the regular case is without recovery, and it follows the EM current density limit which is J, and this will guarantee a 10-year lifetime. If we double this EM current limit and keep the temperature the same, the time to reach each EM stage will be almost halved according to the Black's Equation. But increasing the current density also accelerates the recovery from EM. This paradigm can lead to potential "EM-free" operations over a long lifetime span. Details of the time for different EM stress and recovery phases under various conditions are summarized in Table 3.5. Relaxed EM current rule will lead to less power straps during power network synthesize or less metal layers, and this can potentially save the fabrication cost by reducing the metal stack. It also offers the designer an additional flexible knob and can also bring the performance benefit and less necessary margin, details are discussed in the next sections.

³This condition corresponds to the accelerated test condition in Figure 3.8.

⁴Under this condition, the EM current limit is doubled.



Fig. 3.11 Illustration of cases with and without recovery during normal operation: If the current limit is doubled, the metal wire ages faster (almost twice compared to before), but the in-time accelerated and active recovery can always bring it back to fresh state. Overall, the EM current limit can be relaxed while assuring the reliable operations during the lifetime span.

3.7.2 Performance Improvement

The second benefit that is enabled by EM recovery is the performance improvement due to less IR drop from power and ground mesh. As illustrated in Figure 3.12, when the resistance of the power/ground mesh (ΔR_{PG}) increases due to EM, the performance of the load degrades. To achieve the same performance as before, margins need to be added, but the power consumed by the mesh increases to match the drive current. However, the introduction of EM recovery is able to fix EM and mitigate the resistance increase as shown in Figure 3.11 so that the load always runs at a relatively high speed, the impact of EM-induced IR drop becomes minimal. The potential margins and power overhead that are wasted for matching the performance can be minimized as well. As the changing load can lead to unpredictable dynamic stress conditions, and different metal layer presents different EM



Fig. 3.12 Illustration of IR drop on power mesh.

behaviors [163], the proactively scheduled recovery periods can potentially be a "economic" candidate solution for improving the EM-induced performance loss.

3.7.3 Extend the Wire Lifetime

There are two ways of scheduling EM recovery proactively, the first one is to insert recovery period even during the stress accumulation periods when resistance hasn't changed; The second way is to recover after resistance increase has been detected and reach some level (but are still reversible) just as illustrated in Figure 3.11. Both solutions can lead to that the metal wires are "refreshed" after each stress and recovery cycle. The lifetime of metal lines might not be a bottleneck any more for systems that are expected to operate for a long time reliably. With the recovery implemented on chip, the designers can still design and signoff the chip in a traditional way, but it doesn't need to be lifetime-driven, which means the lifetime specifications for metal wire becomes a reference instead of a constraint. Since recovery is not free as will be discussed in following chapters, the new tradeoff will become how many recovery periods the system could afford vs. the overall lifetime/reliability budgets. The more and earlier the recovery is scheduled, the longer the metal line could last. This offers flexibility of "controlling" the lifetime on the run. For applications that

are lifetime-bonded (e.g. automotive systems, medical devices and implantables), recovery offers opportunities of running for a much longer with an extended lifetime for on-chip metal layers.

3.8 Summary: EM vs. BTI

Based on the experimental results and explored recovery behaviors presented in this Chapter and Chapter 2, we summarize the similarities and differences of EM and BTI behaviors in both stress and recovery periods in Table 3.6. During stress phase, both wearout mechanisms can be accelerated by high temperature under high voltage/current stress. BTI is caused by voltage stress, and the device degradation is characterized by gradual threshold voltage (V_{th}) increase. While EM wearout happens due to the current flow, it increases the metal resistance R, and the degradation behavior is very different from BTI. During early EM stress period, the resistance stays almost the same, and this holds for a relatively long time until the void nucleation phase starts, the resistance presents a sudden increase and eventually reaches the failure state when the wire breaks and no current flows in the metal wire anymore.

In recovery phase, we have demonstrated that high temperature and reverse stress can accelerate and activate the recovery for both phenomenon. While even under these accelerated self-healing conditions, there are still irreversible components which can't be recovered within a reasonable period. We explored that both BTI and EM presents frequency dependency behaviors in which those irreversible parts can be eliminated and fully avoided by scheduling the in-time proactive accelerated self-healing periods. Note that BTI recovery happens only when device is in OFF state, while EM recovery can be activated when devices are ON since there are still reverse current flowing through the metal. If these current can be delivered to the functional block in a proper way, EM recovery will not conflict with the device operating time. Some potential implementations that can enable this will be presented

in Chapter 4. One uniqueness of EM active recovery is that the reverse current can potentially lead to EM wearout in reverse direction. This has been demonstrated by our experimental results shown in Figure 3.8, where extended recovery period after full recovery causes more resistance increase. Thus this suggests that there is an upper limit of EM recovery time, and overly-healing can result in undesired new EM issues. Because of these, EM employs a more complex "circadian rhythm" behavior compared to BTI. Careful considerations need to be taken when scheduling the accelerated and active recovery during implemention, which will be discussed in details in Chapter 5.

Phase	Similarities	Differences		
		BTI - Voltage Stress; EM - Current Stress		
Stross	High T accelerates both	BTI - V_{th} increases gradually; EM - R doesn't		
(Active)	High stress leads to more wearout	increase until void growth and nucleation		
(Active)	ringh stress leads to more wearout	BTI - ΔV_{th} will saturate; EM - <i>R</i> increases		
		(infinitely) until breakdown		
		BTI Active Recovery - Negative Voltage;		
		EM Active Recovery - Reverse Current		
Deservour	High T appalarates both recovery	BTI Recovery - Device OFF;		
		EM Recovery - Device ON (Current flows)		
	Poverse stress activates both recovery	BTI Recovery saturates;		
(Sloop)	Irreversible part can be avoided	EM Over-Recovery can lead to more		
(Sleep)	Inteversible part can be avoided	wearout in the reverse direction		
	requency dependency behaviors	EM has a more complicated		
		"circadian rhythm" due to		
		complex stress behaviors		

Table 3.6 EM vs. BTI - Similarities and Differences during Stress and Recovery

3.9 Conclusions

In this chapter, we mainly focused on recovery techniques for Electromigration wearout that occurs to on-chip metal wires, especially to the power delivery network. Figure 3.13 highlights the main contributions of this chapter. We demonstrated with extensive accelerated tests that EM recovery can be activated by reversing the current direction, and can be


Fig. 3.13 Chapter 3 highlights.

accelerated by raising the temperature. We showed several advantages that can be potentially offered by recovery with an example (in 28nm FDSOI technology node). As EM effects can be rejuvenated by proactive recovery periods, the EM design rules can be relaxed while still guaranteeing reliable operations. The reduced EM-induced IR drop can lead to less performance loss. For system requires a long lifetime, the EM recovery techniques discussed in this chapter can offer such opportunity in an economic way.

The work discussed in this chapter has been published and presented in [C7], [C10], [C11] and [T1].

3.10 Acknowledgements

This part of work was funded by NSF CCF-1255907, SRC Global Research Collaboration (GRC) Program (Task ID. 2410.001), and C-FAR, one of six SRC STARnet Centers, sponsored by MARCO and DARPA. I would like to thank Dr. Linqiang Luo (University of Virginia) for helping with the wire bonding process, and thank Dr. Runjie Zhang (University of Virginia) for discussions during early idea development phase.

Chapter 4

Circuit Techniques for Accelerated and Active Recovery

4.1 Overview

In Chapter 2 and 3, we mainly presented the experimental demonstrations of accelerated and active recovery techniques for BTI and EM respectively. We have also shown that these recovery behaviors were unique and can potential provide significant benefits for future resilient digital system design. As these techniques might not come for free, it is necessary to explore how to utilize the recovery behaviors on chip in an efficient way. This chapter aims to answer the following research questions:

- How to take advantage of the accelerated self-healing techniques on-chip for both BTI and EM wearout?
- What are the potential PPA overhead introduced by the recovery solutions on-chip?
- What are the circuit-level components (such as sensors) and assist circuitry necessary for enabling a true accelerated self-healing system?



Fig. 4.1 Overview of Chapter 4.

Figure 4.1 previews this chapter. As BTI and EM occur at different parts of the circuit, we have introduced separate solutions for each of the wearout mechanisms, this chapter will begin by discussing the circuit implementation for accelerating and activating BTI recovery. A charge-pump based on-chip negative voltage generator is designed and simulated in a 28nm FD-SOI technology. The area of the generator is estimated to be about only $300um^2$, and the ripple is only 1.45%. Since BTI recovery occurs during "sleep", during which leakage is considered a big concern. Power gating techniques [177] have been widely used to "shut off" the leakage paths from logic blocks by inserting power switches (also called sleep transistors) as headers or footers. On top of the existing power gating infrastructures, BTI active recovery can be enabled by adding some extra logic so that the overall overhead is minimal. Since high temperature can accelerate recovery of both BTI and EM, a tiny on-chip heater design that is inspired by [228] will be discussed. The heat generators are enabled only during the recovery intervals, thus the overall power overhead can be leveraged with reliability gains. For the EM accelerated and active recovery implementations, we introduce a multi-mode assist circuit scheme that is able to reverse the current directions while not affecting the functionality of the load blocks. The circuit scheme also supports BTI active recovery mode, during which transistors are reverse biased.

As important components that translate device-level wearout effects into metric degradations that higher levels can understand, on-chip wearout sensors act as monitors for both conventional adaptive solutions and the recovery solutions as discussed in this thesis. The accuracy and reliability of these sensors will be crucial for the system level management units [180] to use their information to monitor both wearout and recovery at the circuit level. This chapter will present three different designs of wearout sensors. Two of them are for tracking BTI wearout and recovery, and one is for sensing EM wearout and recovery. For BTI sensors, the first one is ring oscillator based, where it can continuously track both NBTI or PBTI; the second one is a small metastable element based BTI sensor and serves as a trigger/alert for BTI stress and recovery. These sensors can be potentially deployed across the chip and be reused in various systems to assist the proactively scheduled recovery. Design details, simulation results and physical implementations of all these circuits are going to be detailed in this chapter.

4.2 Circuit Solutions for Activating and Accelerating BTI Recovery

4.2.1 On-Chip Negative Voltage Generation

In Chapter 2.7.1, we demonstrated that negative voltage can significantly boost BTI recovery, but usually there is no negative voltage domain available on-chip, a negative voltage generator is thus needed. In this section, a switch-capacitor (SC) charge-pump negative voltage generator that is modified based on [50] is designed and simulated in a 28nm industry-standard FD-SOI technology node. Shown on the right half of Figure 4.2 is the schematic of the generator. It works as following, during the first half of the charge-pump cycle, *clk*1 is "0" and *clk*2 is "1", the flying capacitor *C*1 is charged to V_{dd} . In the second half, *clk*1 is "1" and *clk*2 is "0", this will lead to that the positive terminal of *C*1 is connected



Fig. 4.2 A switch-cap based negative voltage generator (designed in 28nm FD-SOI technology) for delivering the negative voltage for BTI active recovery. *Vout* is the negative voltage output.



Fig. 4.3 A switch-cap based negative voltage generator (designed in 28nm FD-SOI technology) for delivering the negative voltage for BTI active recovery. *Vout* is the negative voltage output.



Fig. 4.4 Demonstration of integrating the negative voltage generator as part of the PMU (in 130nm bulk technology). The total area of the generator $(300 \times 315 um^2)$ is only $\sim 5\%$ of the total PMU area.

to ground and the negative terminal to *Vout*, *C*1 is in parallel with *C*2. The charge will be redistributed until it reaches the new balance. The value of V_{out} is mainly determined by the capacitance value of *C*1 and *C*2 (*C*1/*C*2) and the clock frequency *f*.

Shown in Figure 4.3 is the simulated results. It shows that the generator outputs a stable voltage of -300.6mV after a startup time of 638ns under a clock frequency of 66.7MHz. One of the most efficient ways of reducing ripples is to use non-overlapping clocking which is also implemented in this work and shown on the left half of Figure 4.2. The resulted ripple is less than 1.45%, which is good enough for our BTI active recovery purposes. The overall area of the generator is only $\sim 4300um^2$ in 28nm technology. The total power consumption is about 64.47uW. While this can be improved by slowing down clock frequency and sizing the switches. Since there are only 4 transistors in the generator, the introduced leakage power overhead is only 68.85nW. The negative voltage generator can be enabled and disabled

by employing the clock gating strategy. During most of the system operating time, the generator is in OFF state unless it is in BTI active recovery period. It is worth to mention that the negative generator has been successfully embedded in a multi-output on-chip switch-capacitor DC-DC converter as part of the power management unit (PMU) for voltage stacking applications [J5] as shown in Figure 4.4, the area of the negative voltage generator design takes only 5% of the total PMU area. In summary, the negative voltage is able to provide stable output voltage and introduces minimal area overhead compared to a regular PMU block in a System-on-chip (SoC).

4.2.2 Negative Bias Voltage in a Logic Path



Fig. 4.5 A chain of inverter logic (FO4) with negative voltage supply: When the voltage supply is negative, and input of the first stage of the logic is "0", the gate and source voltage V_{gs} of PMOS transistors inside each logic stage is positive for all stages of logic so that BTI recovery can be activated. "*nx*" refers to the node voltage at node *x*, *n*0 is equal to 0*V*.

As we have already discussed how to generate the negative voltage supply for activating the BTI recovery in the last section, we now discuss the feasibility of utilizing this generated negative voltage $(-V_{rec})$ in real circuit blocks, which usually includes a chain of logic functions. The question that to be answered in this section is if the negative voltage supply is able to rejuvenate all logic stages in the chain, and if it will introduce overhead such as leakage due to reverse bias across the transistors.

Shown in Figure 4.5 is a chain of inverter logic with Fanout of 4 that is used to mimic the actual logic path in a function block. In this example, the logic depth is only 6 for simplifying the analysis, but in a real IP unit, it can be hundreds or more [195]. During



Fig. 4.6 An equivalent circuit for ease of analyzing the behaviors under active BTI recovery enabled by negative voltage supply $-V_{rec}$.



Fig. 4.7 Simulated internal node voltages of a chain of inverter logic during BTI active recovery with negative supply voltage. X-axis is the applied voltage levels. As the logic depth increases, the node voltage level will increase slowly (an example of recovery voltage at -0.3V shows this).

BTI active recovery period, a negative voltage of $-V_{rec}$ is supplied to the logic path, input V_{in} is set as "0", so PMOS transistors in each logic state experience a positive V_{gs} and BTI recovery can be activated, but not all transistors are biased at the same voltages as $-V_{rec}$. In fact, the intermediate nodes (e.g. n1, n2, etc.) have different values due to the following reasons. shown on the left of Figure 4.6 is an inverter under negative voltage supply, an equivalent circuit is shown on the right half of figure for ease of analysis, it is an inverter in which PMOS and NMOS transistors are flipped with an input equals to V_{dd} . The NMOS turns ON and the output node charges towards V_{rec} . But a delta voltage of ΔV is needed to keep the NMOS in ON state, in super-threshold voltage region, this voltage is roughly V_{th} ; in sub-threshold voltage region, this delta voltage turns to be lower than the threshold voltage. The functionality of this uncommon logic is a weak buffer that passes weak "1" and "0". From the above analysis, the equivalent node voltage for the circuit on the left of Figure 4.6 is equal to $-\Delta V$. As the logic propagate deeper, this delta voltage will increase, thus the recovery voltage across the PMOS transistor will potentially reduce. To investigate what this delta voltage is and how it affects the following logic, we simulate the same circuit structure as shown in 4.2 in a 28nm FD-SOI technology. Figure 4.7 presents the simulated results. The intermediate node voltages nx in each logic stage under different negative recovery voltage level (from $-V_{dd}$ to 0) is shown in the figure. The results suggest that as logic depth increases, the internal node voltage level increases under all recovery voltages. As an example, when $-V_{rec} = -0.3V$, the node voltage drops from -0.0178V (closer to "0") to -0.0435V, this translates directly to the reduced bias voltage V_{gs} across PMOS transistors, and it is shown in Figure 4.8. As we expect larger V_{gs} for each logic stage, the results show the reduction of this voltage, which seems "unfriendly" to the active recovery. But the good news is that this voltage doesn't decrease linearly, it reduces very slowly and saturates in fact as shown in the figure. As a first-order estimation, it takes more than a hundred stage to reach "0" with a recovery voltage of -0.3V. In summary, although applying a negative voltage supply



Fig. 4.8 PMOS Transistor V_{gs} for each logic stage under different negative voltage supply levels. The larger V_{gs} is, the more negative bias PMOS transistors experience, the higher BTI recovery rate will be. As the logic goes deeper, V_{gs} decreases slowly. Depending on the logic depth, negative voltage level needs to be carefully picked.

to a whole function block can't lead to that all PMOS transistors are under the same bias, they experience the positive V_{gs} on different levels, BTI recovery can still be activated as we have demonstrated a very small positive V_{gs} can lead to a big recovery rate boost. Also, to compensate this, the negative voltage can be delivered in a relatively "fine-grained" way so that the PMOS transistors have a larger gate source bias voltage. A second solution is to adjust the negative voltage levels depending on the logic depth so that longer logic chain is able to be rejuvenated as fast as shorter chains.

Another important aspect that needs to be carefully evaluated when introducing a negative voltage supply in the logic is the power consumption. As the active recovery occurs when transistors are OFF, thus leakage power is a main metric of interest. As shown in Figure 4.6, when under a negative voltage supply, the PMOS transistor is in forward bias (FBB)



Fig. 4.9 Leakage power consumption during BTI active recovery. All values are normalized to the case when logic is OFF without any negative bias ($V_{dd} = nominal \ voltage, V_{in} = 0 \ or \ V_{dd}$).

condition, under which V_{bs} is less than 0, and this can lead to a decreased V_{th} and increased leakage current $I_{leakage}$ potentially. To explore this, simulation with the same setup as in Figure 4.7 is conducted and results are presented in Figure 4.9. As a baseline case for comparison, the leakage power during passive recovery condition (V_{dd} is the nominal voltage, V_{in} is constant "1" or "0") is used for normalization. All leakage power values ($I_{leakage} \times V_{rec}$) presented in the figure are normalized value. The results surprisingly show that the negative voltage doesn't introduce any leakage power overhead even when the negative voltage is biased at full range (-1V in this case). Since we demonstrated with experiments that only a small negative voltage (such as -0.3V) is good enough for boosting the recovery, the leakage power can actually be reduced significantly. Thus we conclude that it is feasible to implement



Fig. 4.10 Power gating structure: it can be used effectively to reduce leakage and help BTI recover passively, retention registers are used for saving the states. The voltage drop due to the resistance of header transistor ΔR can lead to performance loss, it is even worse that this transistor is ON most of the time and experience BTI wearout as well.

a negative voltage supply on-chip, and the solution is able to efficiently rejuvenate a circuit block from BTI without introducing additional leakage power overhead.

4.2.3 Wearout-aware Power Gating

Power-gating was firstly introduced as an effective method to reduce the standby leakage by inserting sleep transistors (ST) between the logic blocks and the actual power/ground rails [36]. Figure 4.10 shows a version of power gating with header sleep transistor inserted. Retention registers are used for storing the states when logic is in sleep mode. Several work [34–36, 132] have shown that the insertion of the sleep transistors also makes the circuit immune to BTI-induced wearout due to that the structure enables the passive recovery mode by generating idle periods. While the implementation of power-gating doesn't come for free. Besides the area overhead, the sleep transistors also introduce the on-resistance ΔR , which



Fig. 4.11 Sleep transistor threshold voltage increase ΔV_{th} vs. Load performance loss (8 49-Ring Oscillators running in parallel, the header transistor is sized as $1\mu m$ wide).

further leads to performance penalty. The size of the sleep transistors, together with the size of the power-gated blocks also determine the power down or up time [36]. When BTI wearout is introduced into the picture, power gating design becomes more challenging. On one hand, power-gating is preferred for all the logic, especially the critical path, so that it could be recovered during idle time, on another hand, the introduced delay overhead is not preferred. A clustered architecture was proposed in [34, 36] to leverage power, performance and lifetime tradeoff by using two types of sleep transistors. Also in [152], a microarchitecture-level framework was proposed to mitigate wearout by utilizing power-gating as a design knob to mitigate the wearout effect for superscalar processors. A numerical model was proposed in [39] to analyze the potential benefits and limitations of power gating for reducing BTI. While one fact that is ignored in these previous work is that sleep transistors are almost ON most of the time, so they experienced wearout the same way as the gated logic blocks. As a result, the circuit performance and lifetime become even worse. This is shown in Figure 4.11, where we simulate the threshold voltage increase (ΔV_{th}) of sleep transistor vs. load performance loss in 28nm FD-SOI technology. The load are eight 49-stage ring oscillators running in



Fig. 4.12 An active recovery-enabled power gating structure: Blue outline parts refer to the additional logic on top of the existing infrastructure. Sleep signal is the trigger for starting the active recovery.

parallel. This performance degradation can be as large as more than 10% under a small threshold voltage shift. To tackle this problem, [34] proposed to realize NBTI-aware power gating by oversizing the sleep transistors, using forward bias to compensate and reducing the stress time. [119] explored the sleep transistor redundancy to reduce overall turned-on times of these transistors. [232] further analyzed the joint interdependent degradation effects on logic networks and sleep transistors by using redundant STs. In [233], two BTI-aware sleep transistor sizing algorithms were proposed to reduce the total width of sleep transistors based on the distributed sleep transistor network structure. [176] performed an analysis and design flow to pick threshold voltage values for STs to optimize both leakage and lifetime.

On top of these previous solutions, this section discusses an orthogonal power gating solution that enables full BTI recovery by delivering the negative voltage supply to logic blocks. The structure also incorporates the rejuvenation for the sleep transistor by using a higher-than- V_{dd} voltage, and it is shown in Figure 4.12. An NMOS switch is added to the virtual power supply node for delivering the negative voltage when active recovery mode is



Fig. 4.13 Functional simulation of the proposed wearout-aware power gating structure in 14nm bulk FinFET technology, signal names correspond to the ones in Figure 4.12: When *Sleep* signal is high, the BTI active recovery mode starts, the negative voltage is delivered as the virtual supply for the logic, the header sleep transistor is also in reverse biased because a higher-than- V_{dd} voltage is applied to the gate. The switching time between modes are similar to the ones in existing power gating infrastructure.

activated (Signal *Sleep* = 1). Since the capacitor *C*2 in the negative voltage generator (shown in Figure 4.2) acts as the decoupling capacitor (Decap) explicitly, so no external Decap for the negative voltage is needed. Figure 4.13 is the functional simulation result of the proposed power gating structure in an industry-standard 14nm FinFET technology node. The load are four 9-stage ring oscillators running in parallel. The header transistor is sized so that the voltage drop is within 5%. V_{dd} _high is selected so that the negative bias across the PMOS header is $\sim -0.3V$. The NMOS transistor is added to deliver the negative voltage of -0.3Vfrom the voltage generator. It is sized so that the switching time between modes are similar to the regular power gating case where there is no active recovery implementations. As shown in the figure, when BTI active recovery mode starts, the structure is able to supply a very stable negative voltage to the load. When the *Sleep* signal is set as "low", it can switch to



Fig. 4.14 The physical layout of the power gating that implements the BTI active recovery logic in 28nm FD-SOI technology. The added overhead is only a NMOS header and some control logics. The load is 8 ring oscillators in parallel, and the frequency divider is used for off-chip frequency readout.

the normal operating mode from recovery mode quickly. This switching time can be further reduced by leveraging the load size and the NMOS header transistor size.

The proposed wearout-aware power gating structure introduces some extra logic such as a buffer and a NMOS transistor, but the overhead of these components are very minimal. For example, in the case shown in Figure 4.13, the added area overhead is only $2.71um^2$, which is close to area of a X5 Flip-flop in the same technology. This can be seen in Figure 4.14, where it gives a physical implementation of the wearout-aware design. We integrate the BTI active recovery with the existing power gating structures. The load are 8 ring oscillators in parallel, and the frequency divider is used for offchip frequency readeout. The added logic includes a NMOS header (shown on the top left part) and control logic, which includes an input buffer and output MUXes. The main area overhead can come from routing as two more voltage domains (higher-than- V_{dd} and negative voltage) are introduced, delivering these voltages to logic can take some routing resources, but the good news is that most of the modern power gating structures are employed in a relatively coarse-grain way, such as at core level or a functional block level. So we expect that this routing overhead can be improved by a careful designed floorplan and optimized routing options. As for the source of the higher-than- V_{dd} voltage, designing a voltage regulator for generating it can be very costly, but two alternatives can be implemented. The first one is to combine with the body bias techniques (e.g. [93, 94]) and utilize the generated high voltage for recovery purposes; A second potential solution

is to utilize the voltage from other power domains. Since most of the modern SoCs have multiple power domains for improving the energy efficiency [100], this offers a possible way that delivers the higher voltage for recovering the logic for lower voltage domain. Details of the implementations depend on the specific applications and design specifications. As a summary, the proposed "recoverable" power gating design is orthogonal to the existing power gating structures, it can be easily integrated with these infrastructures and optimization framework to achieve the ultimate goals of ultra-robust digital systems.

4.2.4 **On-Chip Heat Generation**

As demonstrated in Chapter 2 and 3, high temperature is able to accelerate both BTI and EM recovery significantly, but it is also well known that heat can degrade the performance. To this contradiction, we propose to use small reconfigurable localized on-chip heaters for accelerated recovery in this section. These heater designs are inspired by previous work [10, 11, 228], where they were used as self-heating elements for thermal-aware testing on FPGAs. Thus, similar designs can be adapted to our high temperature-enabled accelerated recovery purposes. A proof-of-concept design is presented in Figure 4.15. The theory behind it is to force high toggling rates of the a reconfigurable ring oscillator to generate heat, and the output frequency can be selected based on the temperature requirement. The design is straightforward but easy to implement and control. *Accelerated Recovery* signal serves as an enable signal and can be triggered by recovery decisions. Note that the length doesn't need to be always half of the total length, it depends on temperature requirement that is necessary for recovery.

Since each frequency corresponds to different temperatures, it is important to understand how these two metrics are related. Figure 4.16 illustrates a case where similar heating designs are implemented on FPGAs, and the measured maximum temperatures during steady state are plotted [228]. Frequency and temperature have an almost linear relationship, where



Fig. 4.15 Reconfigurable On-chip Heating Elements Schematic: The output frequency of the ring oscillators can be reconfigured by selecting the length of the inverter chain (L, L/2, etc.). Different frequencies of the heating elements correspond to a wide range of temperatures.

temperature increases by 10°C with the increase of $100MH_z$. While it should be noted that heaters need to take some time to reach a steady state due to thermal resistance. This time can be in the order of minutes, which we believe are still acceptable for recovery purpose. Authors in that work also observed that temperature swing with respect to idle temperature increases with higher frequency of the generator, and this swing can be as large as $\pm 5^{\circ}$ C. But our experiments results show that this amount temperature swing in high temperature range (> 80°C) has a relatively minor impact on recovery rate. Although the relationship in Figure



Fig. 4.16 Maximum temperatures that correspond to different oscillation frequencies with the heating elements on FPGAs, temperatures are sampled with external thermal sensors, a precalibration has to be carried out on the target FPGA to determine this relation, data is from [228].

4.16 might not be always the same for other designs or chips due to the differences of thermal resistance, it still indicates that the on-chip heater is able to achieve a wide range of relatively steady high temperatures. A precalibration process is usually necessary to determine the exact relationship for the target devices [11].

As the distribution of the generated heat from the heater across the desired region is totally controlled by the placement of these heating elements. A homogeneous temperature distribution across the block is desired for recovering all the logics in that block. Here we discuss three different strategies of placing heater on the chip. These are shown in Figure 4.17. The first one is the most ideal case where the heating elements are spread evenly inside a logic blocks so that a homogeneous temperature distribution can be achieved across the block area. In fact, this spreading can lead to that the temperature on the edges of the block is



Fig. 4.17 Potential strategies of on-chip reconfigurable heaters: (a) Evenly distribution; (b) Ring placement; (c) Critical path placement. Red squares refer to the heating element. An external controller is needed to configure the heater and select the blocks that need to be recovered.

slightly lower than at the center. To handle this, the heater in the center can be reconfigured to oscillate with a slower frequency and thus balance the temperature across the whole area. The difficulties of implementing this strategy is that it requires much more design efforts such as placement and route, especially for logic blocks that are not rectangular or square. A relatively easier way to place the heating elements is shown in (b) of Figure 4.17, where a "heater ring" is formed and placed surrounding the logic block. The advantage of this method is that the heating elements are separated from the logic, so they have less effect on performance during the run time. But the disadvantage is that the center of the logic block is not able to reach the desired recovery temperature due to the existing of physical distances. Thus this solution is applicable to those smaller blocks or those logic that are less wearout critical, and this heater ring can be potentially be shared by logic blocks that are next to each other. As the required number of heating elements are proportional to the logic area with this ring placement method, another more economical way of placing them is shown in (c) of Figure 4.17. The heaters are placed only close to critical paths which determine the overall

performance of a logic block. The number of heaters can be reduced, and it can be potentially more effective. The challenge is to have the right methodology to place them during the early design phase. One possible solution could be to embed the single element into a scan chain cell and place them during the DFT process. As many heating elements are expected on a single chip, a heating controller is required to configure the heaters individually, it also needs to activate/deactivate a subset of heating elements according to the desired maximum temperature. This controller can be a look-up-table (LUT), which takes the desired recover temperature from the user as the input, and compares it to the FPGA current temperature, which can be either read from the built-in thermal sensor available on chip, or predefined by thermal analysis tools, and then output the corresponding select signals for the right group of heaters. In summary, each of the above three placement strategies has some advantages and disadvantages. Careful design decisions that take advantages of the high-temperature enabled recovery and balance all tradeoffs can lead to effective and complete self-healing.

As shown in Figure 4.15, each heat element composes a chain of inverters and a MUX. The length of the chain depends on the desired temperature, and thus the desired frequency. As an example, for an oscillating frequency of more than *GHz* range (corresponding to $> 80^{\circ}$ C), about ~ 41 stages of X4 inverters are required with a 28nm FD-SOI technology, the total area is around $16\mu m^2$, and the leakage is about 16.8nW. This is still much less than a regular on-chip temperature sensor (such as [234, 224]). Thus the introduced additional area and leakage are acceptable. But it should be noted that the on-chip heating elements are based on the assumptions that extra power consumption are allowed during recovery. For example, this power can be as large as mW depending on the target temperatures. It might not be suitable for all applications, especially for those have very tight energy efficiency requirements. But the tradeoffs are between the energy and power consumption during OFF period and the improved metrics during active time. Careful designs can potentially leverage these two so that the overall overhead is minimal and acceptable. An alternative

and potentially more economical way of recovering with temperature is to utilize core level redundancy and have active cores/elements serve as "natural on-chip heaters" for the neighboring asleep cores. This will be further discussed in details in Chapter 5.

4.3 Circuit Solutions for Activating and Accelerating EM Recovery

The experiments in Chapter 4 demonstrated that not only high temperature is able to rejuvenate the on-chip metal lines from EM damage, revering the current direction is also an effective recovery boost solution. Since power rails suffer from single-direction DC current mostly [86, 204], we focus only on EM-induced effects in power delivery network in this thesis. To enable the on-chip implementation of EM active recovery, an assist circuit scheme that is able to deliver the reverse current when needed without affecting the functionality of the logic is presented in this section. The idea is inspired by the concept proposed in [4, 16, 185]; the difference between this scheme and previous solutions is that our scheme is able to support both EM and BTI active recovery modes, and we also discuss physical implementations and potential system level integration, which are missing in the literature. Shown in Figure 4.18a is the schematic of the assist circuitry, which includes two layers of power gating, with power grids in between. Functionally, it supports three modes (Normal, EM Active Recovery and BTI Active Recovery). Figure 4.18b is the corresponding truth table. A close look of these three modes are is presented in Figure 4.19. Under Normal operating mode, the functional load works similarly to a regular power-gated system, and during EM Active Recovery, the current flowing through the VDD and VSS grid is reversed, and the current has the same absolute value that is guaranteed by the symmetry of the scheme, thus the load (target circuit) still functions as under Normal mode. BTI active recovery happens when the load is idle, during which VDD and VSS of the load are switched. This



Fig. 4.18 Assist circuitry for activating BTI and EM recovery: (a) The main circuitry, arrows represent the current direction under different modes, V_{DD} and V_{SS} pins can be connected to the on-chip voltage regulator directly, or to the global power delivery network; (b) Truth table for three operating modes; (c) An example of activating NBTI recovery under *BTI Active Recovery* mode, for PBTI recovery, the input needs to be "0", ΔV represents voltage droop/increase or noise.

creates a BTI active recovery condition not unlike what has been described in Section 4.2.2. Depending on the input values, NBTI or PBTI recovery can be activated; this is shown in Figure 4.18c.

To validate the design, we implement and simulate the above recovery assist circuitry in a 28nm FD-SOI technology. A set of ring oscillators running in parallel is used as the load, the



Fig. 4.19 Three modes of the assist circuitry for activating BTI and EM recovery: *Normal Operation*, *EM Active Recovery* and *BTI Active Recovery*. VDD and VSS grid in the figure refer to power delivery network (PDN).

VDD/VSS grid is treated as a resistor for which we picked a reasonable value based on the published literature. Figure 4.20 presents the functionality simulation under three different modes. Figure 4.20a is the current flowing across the VDD Grid during under *Normal* mode and *EM Active Recovery Mode*, the direction of the current is reversed, but the amplitude of two currents are still the same, this ensures that the load can still run with the same



(a) The current direction is reversed under *EM Active Recovery* Mode, but the current value is still the same so that the load still runs at the same frequency.



(b) Under *BTI Active Recovery* Mode, load VDD and VSS values are switched so that the voltage across the transistors are reverse biased.

Fig. 4.20 Functionality Simulation for EM/BTI Active Recovery Assist Circuitry in 28nm FD-SOI technology.



Fig. 4.21 Load Size vs. Performance and Switching Time: Increasing the number of loads will reduce the performance and increase the switching time between modes, to compensate the degradation, header/footer transistors need to be upsized, which will further increase the area. This tradeoff needs to be carefully considered during the design process.

performance. Under the *BTI Active Recovery* mode, the load is in sleep mode, its VDD and VSS nodes are switched as expected, and there is about 0.2V voltage droop/increase induced by the pass transistors, but the voltage is still large enough for activating BTI recovery (-0.816V is much higher than -0.3V demonstrated by our experiments in Chapter 2). Three modes are configured by a 3-bit controller, which is a decoder following the truth table in Figure 4.18b.

One of the biggest design challenges of the active recovery assist circuitry is the voltage droop/increase at the load VDD/VSS nodes that are introduced by two layers of header or footer transistors during *Normal* operation and *EM active recovery* mode, during which load performance is critical. Another potential concern is the switching time (retention time) between modes. Since both metrics depend on the load, we explore how load size affects them. Figure 4.21 shows that by increasing load size, the performance degrades linearly



Fig. 4.22 Vertical cross section of the physical implementation for the EM/BTI Active Recovery Assist Circuitry for VDD Grid (VSS Grid can be implemented in a similar way): EM hazards happen at high current density regions, which could be caused by faster switching activities on the load logic; At the logic level, BTI hazards happen due to the continuous stress.

because of the voltage drop/increase across the footer/header transistors. Switching time also increases with the increased load, but with a slower rate. To compensate this performance degradation, the header/footer transistors need to be upsized accordingly, which will results in more area. This study indicates that each load will have its own optimal design point which gives the optimal values in terms of area and other metrics. Thus, careful design decisions are necessary to balance the tradeoffs of performance vs. area vs. reliability.

Figure 4.22 illustrates an example of physical implementation of VDD Grid with the assist circuitry integrated (10-layer metal stack is assumed, and this is what has been used in many advanced technology nodes such as 28nm FD-SOI and 14nm FinFET). It has a global PDN grid which is usually built with the top one or two metals that are wide and

thick, thus being relatively robust against EM. Local VDD/GND grids that are close to logic and use the lower metal layers are more EM sensitive; *P*1 to *P*4 correspond to the power gating header transistors in Figure 4.18a. This implementation will be able to protect the local grids and also enable the flexibility of designing localized assist circuitry for individual loads that are affected different by wearrout. The structure is very similar to a conventional power gated PDN, on top of which we add one more layer of header/footer. Since power gating techniques have been widely used, this implementation makes it easy to integrate the assist circuitry into the existing design flow and existing infrastructures.

The overall area overhead introduced by the assist circuitry are determined by design specifications in terms of maximum IR drop, and also depends on load size and behaviors. The tradeoffs can be leveraged and balanced based on the applications and budgets for power and area costs.

4.4 BTI Sensing

As BTI wearout happens at the transistor level, and some of the recovery decisions or adaptive solutions occur at the circuit or higher level of the system stack, on-chip wearout sensors have been an important component that bridges these two levels. They act as monitors for both conventional adaptive solutions [114, 146] and the recovery solutions as discussed in this thesis. The accuracy and reliability of these sensors will be crucial for the system level management units [180] to use their information to monitor both wearout and *recovery* at the circuit level. We can expect the total number of on-chip wearout sensors in many-core systems to reach thousands [98, 180], with hundreds of sensors per core. There have been various sensor designs recently that use advanced circuit techniques to achieve good sensitivity in the presence of variations. However, many of these are analog in nature, and are also not able to track *recovery*, thus this limits their wide deployment and reuse in a big system. In this section, we present two types of BTI sensing techniques, the first one is Ring

Oscillator-based and can differentiate NBTI and PBTI, thus the structure can be used for BTI testing and charactering. The second one is a novel type of small and embeddable digital BTI wearout sensors that are based on metastable elements. It is able to track both wearout and *recovery*, and is also able to detect BTI-induced critical path reranking [140, 223] by tracking multiple paths simultaneously. This section details the design, implementation and placement methodology of the novel BTI sensors.

4.4.1 Previous BTI Sensing Techniques

The effect of BTI can be mapped to several metrics, the most straightforward one being frequency degradation — this is also the most common way of detecting BTI [106]. In [90, 108, 186], ring oscillator (RO) based sensors were proposed using tracking and reference ROs to detect the frequency difference. These sensors can achieve very fast tracking, but area and power overhead are huge. Additionally, intra-die variations of the ring oscillators will make the measurement inaccurate. Some work also proposed to use on-chip circuit such as SRAM to predict wearout [111]. Another metric that can be used for BTI sensing is leakage current or drain current — there are some work that use highly accurate current testing techniques to track BTI [51, 96], but it is difficult to use these techniques for real time monitoring; also these circuits are usually very complex and need off-chip measurements. In [123], a self-testing technique for tracking NBTI was proposed — this sensing technique is robust across the PVT variations, but it is still ring oscillator based and has huge area overhead. [231] proposed a small embedded NBTI sensor based on a metastable circuit. The circuit consists of a pair of cross-coupled inverters. Instead of providing data at every time point, these sensors can provide feedback for the system when the circuits have degraded by a certain critical percentage ΔI . This schemes allows smaller sensor sizes while still achieving good accuracy. [205] extended the work by designing a fine grained sensor that can track both N/PBTI. The resolution time measurement is done by using a ring oscillator-based

time-to-digital converter. [8] proposed an on-chip NBTI sensor which is able to isolate the NBTI and PBTI aging effects as well. But it requires multiple copies of the same sensor, and this could lead to large area overhead. Machine learning-based BTI sensing techniques appeared recently, the idea was to identify the delay degradation for a large pool of paths through monitoring the delay of a selected small subset of paths [113, 60]. As there is no requirement to equip every critical path with a BTI sensor, the area overhead can be reduced.

through monitoring the delay of a selected small subset of paths [113, 60]. As there is no requirement to equip every critical path with a BTI sensor, the area overhead can be reduced. In summary, all these previous work focused on tracking the degradation of BTI of critical path/block and didn't really consider how to embed these sensors into a big system and/or design flow. Also to the best of knowledge, there is no BTI sensor that has been designed for tracking proactive recovery, and tracking multiple (critical) paths simultaneously with one embedded sensing structure. So if critical path reranking [223, 140] happens, the sensor information will not be accurate any more. In this section, we present two novel types of BTI sensors, in which the first type is able to separate NBTI and PBTI with only one test structure. The second type is able to support both stress and recovery sensing with an ability of tracking multiple (e.g. 4 or more) critical or potential critical paths simultaneously. The compatibility of the sensors makes them easily included as a single IP block with the conventional top-down ASIC design flow.

4.4.2 Ring Oscillator-based Test Structures for Separating NBTI and PBTI

PBTI was used to be considered as a secondary effect compared to NBTI. But it now has grown to be a considerable reliability concern as high-k dielectric material and metal gate are adopted for gate leakage reduction [237]. Sensing or charactering PBTI is as important as for NBTI. Conventional sensing techniques or test structures such as ring oscillators don't allow individual and independent measurement of PBTI and NBTI. Existing work [8, 110] proposed solutions of isolating NBTI and PBTI by including multiple copies of the



Fig. 4.23 Sensing element for separating NBTI and PBTI. (a) Functional schematic of 1 sensing stage. Red refers to the actual logic that experiences N/PBTI. Odd number of stages can form a ring oscillator and serve as test structures or sensors. (b) Truth table for the sensing stage. When under NBTI or PBTI mode, PMOS or NMOS transistors are under constant stress; when under Test mode, the structure works as a regular ring oscillator stage, and the oscillation frequency can be read out from the RO output.



Fig. 4.24 Layout of an example of using sensing element in a 37-stage ring oscillator in 28nm FD-SOI technology. Overall area is small, and the sensor can be configured to sense either NBTI or PBTI. It can also be used as the test structure for studying N/PBTI degradations.

sensing structure and configuring them for different modes. These solutions work effectively, but large area overhead can be introduced due to the fact that multiple sensors might be required for each BTI mechanism in a large SoC. Thus, we present a simple and integrated ring oscillator-based sensing stage that can monitor NBTI and PBTI separately while not introducing significant area overhead. Figure 4.23a shows the schematic of one such stage,

where the red inverter is the circuit of interest which experiences NBTI or PBTI. If the ring oscillator cells are replaced with this structure, the oscillation frequency degradation will be able to reflect either PBTI or NBTI depending on the stress mode. The transmission gate serves as a switch for Stress mode (NBTI or PBTI) and test mode where the oscillation frequency of the ring oscillator is sampled. Figure 4.23b is the truth table for the functionality of the sensing stage. Under NBTI mode, the circuit is configured such that only PMOS is under constant stress (full V_{dd} swing) in each stage. PBTI mode works similarly. When under test mode, the sensing stage is configured as a regular single-stage inverter, and can as a ring oscillator cell. The output frequency of the ring oscillator is sampled under this mode.

The design has been implemented in 28nm FD-SOI technology. Figure 4.24 is the layout for a 37-stage ring oscillator that is implemented with the sensing element discussed above. The overall area is about $92\mu m^2$, and it is only $\sim 1.6 \times$ larger than the same RO configuration that is implemented with single inverter only. The structure can be used as a building block for either NBTI or PBTI sensing scheme, and it can also be employed as a BTI test structure for characterizing N/PBTI and understanding the differences between two mechanisms.

4.4.3 Metastable-Element-based Embeddable BTI Sensors

The critical path determines the performance of the system, but, due to the fact that BTI wearout is a gradual effect that accumulates with time, different paths experience various switching activities and/or temperatures, it is highly possible that some non-critical paths, that are only marginally faster than the critical path become the new critical path (and vice versa) in multi-output cases [223]. Figure 4.25 gives an example of such scenario. The circuit shown in Figure 4.25a is simulated in a 28nm FD-SOI technology node. BTI-related wearout parameters for that technology are extracted based on the experimental results from [53, 14]. Assume *path*1 is (part of) the critical path, and *path*2 is the second slowest path in a circuit. Given different input patterns for the two path as shown in Figure 4.25b,



(c) Path delay after BTI stress. Path reranking happens after around 2.3 years, the delay of path 2 surpasses path 1.

Fig. 4.25 Illustration of the BTI-induced Critical Path Reranking

simulated results in Figure 4.25c shows that after certain times, the delay of *path*2 surpasses *path*1, and this results in so-called path reranking (or reordering). The conventional ways of dealing with BTI wearout and sensing BTI mostly focused on the critical path without considering the path reranking issues. If this is left unchecked, the system will experience an unexpected performance penalty. For BTI sensors, tracking more paths, especially the second or third most critical path is thus necessary to guarantee the expected performance at the system level. The most obvious way of doing this is to use individual sensors for each path, but due to variations (like process) of these sensors, one path usually needs many sensors, the area overhead becomes unacceptable. So compact sensors that are able to track multiple paths simultaneously are preferred in this case. This section will present one such sensor which is able to track BTI recovery so that any proactive recovery techniques discussed in previous chapters can be enabled by the sensor outputs. The proposed sensors are implemented and simulated in an industry standard 28nm FD-SOI technology. Details will be discussed in the following sections.

Sensor Scheme

Figure 4.26 shows the main topology of the proposed NBTI sensors. The sensor is based on the metastable cell which composes of a degradation inverter and a reference inverter [231, 205]. Driven by the paths of interests, the sensor is exposed to the same environment as the core circuit, and degrades the same way as the paths. Figure 4.26a shows the high level block drives several cells for active recovery purpose. Figure 4.26b is the case when multiple paths drive separate cells, and is aware of path reranking. In both cases, only one reference inverter is needed and can be shared by all tracking inverters, and the outputs can also be used as trigger signals for DVFS or body biasing techniques. Figure 4.26a and b into a



Fig. 4.26 Multiple-Critical-Path Embeddable BTI Sensor high-level scheme. (a) Proactive Recovery Sensing; (b) Multiple-Paths Sensing.

single design. The bottom half of the schematic corresponds to 4.26a, where P0 - P2 are tracking transistors that are sized based on the degradation percentage (2%, 5% and 10% in this case). The percentage is determined by the driven strength of the transistor, which further corresponds to the size. Similarly, the top half corresponds to 4.26b which enables multiple paths tracking with the same degradation percentage limit (5%) for detecting the path reranking issues. Each half has a reference transistor (*P*6,*P*7). The binary outputs *a* and *b* are read out through inverters. The transmission gates between tracking transistors and paths are used for selecting certain path, and also reducing the critical path performance impact due to the sensor. The symmetric structure is able to leverage the load imbalance between node *a* and *b*. In this implementation, the sensor is for NBTI tracking, PBTI tracking can be enabled by replacing all the PMOS logic to the corresponding NMOS logic. The design shown in the figure can track 4 paths. More can be extended by copying one branch and adding to the existing one if needed.

Figure 4.28 shows the sensor functionality through simulation. The inputs can be encoded to only 6 signals. The sensor has three modes, tracking mode, polling mode and recovery


Fig. 4.27 Transistor-level schematic of the BTI sensor.



Fig. 4.28 Functionality simulation of the BTI sensor under fresh (time = 0) condition.

mode. In tracking mode, the track signal is set high and tracking transistors are stressed by different paths. The reference transistors are OFF to keep "fresh". NMOS switches S0 - S7 are turned off to separate the tracking transistors from the polling circuitry. In polling mode, the *track* signal is set as low, S6 is turned on if right half paths need to be polled. Similarly, S7 corresponds to the left half paths. S0 - S5 are turned ON one by one, during each ON period (polling period), both outputs (*a* and *b*) are discharged quickly, and then poll signal is set high, a fight condition is created between the tracking transistors and the reference transistors, since tracking ones are sized larger (stronger) initially, assume that the loads seen from each inverter are equal, stronger ones will "win the fights" and latch the outputs to high. When NBTI-induced degradation limit is reached and reaches a point when reference transistors win the fights, it means the expected degradation is clear that it cannot introduce significant NBTI-induced degradations or recovery from a long term perspective. Recovery mode works similar to the tracking mode, r0 - r5 are recovery control signals for each path. When recovery is scheduled, these signals are set to low, and the tracking transistors

are in recovery mode, during these recovery periods, polling is still needed to check if the expected recovery level is achieved. A nice feature of the sensor is that any recovery boosting techniques discussed in Chapter 2 can also be applied to the sensor easily so that it is under exactly same recovery conditions as other circuitry.

Sizing

Size corresponds to the strength of the transistors. Tracking transistors are sized x% larger (*x* is set as the threshold for BTI wearout) than reference transistors using multi-fingers for easy area changes. Since variations will affect the accuracy of the sensors, we pick the finger size (120*nm*) a bit larger than the minimum (80*nm* in 28nm FD-SOI technology). All switches are sized equally. Two footer NMOS transistors (*N*0 and *N*1) in the metastable cell are also sized equally and relatively large to have a fast discharging. The load imbalance is dealt with by adding dummy transistors to the load as suggested in [231]. Variations are also dealt with by using interdigitation techniques for degradation transistors reference transistors in physical design phase.



Fig. 4.29 Simulation setup for introducing BTI wearout in a circuit netlist.

Proactive Recovery Case Simulation Results

All simulations are done with Cadence Spectre. The NBTI-induced threshold voltage shift is introduced into the circuit netlist by adding a DC voltage source at the gate of the

tracking transistors as shown in Figure 4.29, degradation is calculated from the BTI models discussed in Chapter 2. The operating conditions are at 1V and 398K.



Fig. 4.30 Threshold degradation and sensor trigger point.

Proactive recovery corresponds to the bottom half of the sensor as shown in Figure 4.27. Assume *path*0 is always at "low" so that all tracking transistors on the bottom half are under constant stress. Figure 4.30 shows the threshold voltage shift ΔV_{th} during stress and recovery, where flags indicate when the sensor is triggered. Figure 4.31 is the corresponding triggering order of the sensor. We check *outa* in this case. ① corresponds to the fresh status when no tracking branch is triggered. When the threshold voltage shift ΔV_{th} is equal to 9.5*mV*, the 2% branch is triggered (②), and ΔV_{th} keeps increasing until 5% threshold is hit (③, ΔV_{th} is 15.8*mV*). Then the recovery signal is triggered, the recovery starts until it reaches the point when all branches are not triggered (④), and the process will repeat. Instead of waiting until 10% path is triggered (ΔV_{th} is 25.4mV), the proposed sensor is triggered much earlier before the end of life so that proper active recovery techniques can be applied to prevent early-life failures.



Fig. 4.31 Sensor triggering order for the proactive recovery case.

Multiple Critical Paths Case Simulation Results

In this case, the main functional part of the sensor is the paths on the top of Figure 4.27. Assume that *path*1 is the critical path, *path*2 and *path*3 are possible critical paths with different switching activities. Three different input patterns are given in Figure 4.32. Figure 4.33 is the corresponding sensor triggering order. The earliest triggered path (*path*2) shows the highest sensitivity to NBTI. Both *path*2 and *path*3 experience more degradation than *path*1. Thus critical path reranking occurs. With this detected by the BTI sensors, some



Fig. 4.32 Input patterns for multiple-path simulation case.

compensation techniques such as high voltages or back biasing can be enabled by these sensor outputs to avoid the performance degradation in the early lifetime.

Tolerance to Process Variations

To show that the BTI sensors discussed above are robust across process variations, Monte Carlo simulations are run with both intra-die and inter-die variations modeled by the foundry. Figure shows the results when sensor is in "fresh" status. The outputs are always correct across the variations. After degradation, at more than 85% of the points (500 points in total), the sensor outputs are the same as at the TT corner.

Physical Implementations and Comparisons

Shown in Figure 4.35 is the physical implementation of a 2-path version of the proposed BTI sensor in 28nm FD-SOI technology. The overall area is about $54.1um^2$. The leakage



Fig. 4.33 Sensor trigger order in Multiple Critical Paths case.

power is 40.64nW. Table 4.1 compares different metrics of the proposed BTI sensor against other state of art circuit level BTI sensor designs. Compared to the ring oscillator-based sensors, the metastable sensors are smaller and faster in terms of data acquisition time. The type of metastable-element-based sensors presented in this section takes the least area and is



Fig. 4.34 Monte Carlo Simulations at the "fresh" status (500 points). More than 85% of the seeds give the correct outputs. The overlapping color in the figure refer to waveforms for each seed. The output of the sensor at TT corner is used as the reference.



Fig. 4.35 Layout of a 2-path version of the proposed BTI sensor in 28nm FD-SOI technology. The total area is only $54.1um^2$.

able to support multiple path tracking. As they are tiny, they consume least leakage power and are fast in sensor output acquisition. Thus this type of sensor can be potentially deployed and distributed in a big design such as multicore systems. A number of sensors can be spread over each core to achieve statistical results. Depending on area and reliability requirements, the number of paths tracked by the sensor can be picked correspondingly. The outputs of the sensors can trigger the conventional adaptive techniques, like Dynamic Voltage Frequency Scaling (DVFS), adaptive body biasing and voltage scaling methods to compensate the loss.

Туре	Work	Tech.	Area (μm^2)	Leakage	Function	Acq. Time	# of paths
Ring	[102]	130nm	277950	-	Wearout	29us	-
Oscillator	[103]	45nm	150	-	Wearout	-	1
Based	[194]	45nm	77.3	-	Wearout	-	1
Metastable	[231]	150nm	493.2	-	Wearout	200 <i>ns</i>	1
Element	[205]	32nm	105	239nW	Wearout	-	1
Based	This	28nm	54.1	10 61 mW	Wearout/	1225	2
	1 1115	201111	54.1	+0.04////	Recovery	12/13	

Table 4.1 Comparisons against other BTI sensor designs (all are circuit-level BTI sensors)

Orthogonally, the ability of tracking recovery makes the sensor a perfect trigger for the proactive accelerated and active recovery techniques discussed in this thesis.

Sensor Placement Methodology



Fig. 4.36 A 2-path version of the sensor implemented with the PnR tool directly.

Including the proposed sensor into a top-down ASIC design flow will be more efficient than the custom design solution, and it also enables the design reuse, especially when a large number of sensors are distributed in a complex circuit. As all instances (single transistors, transmission gates or inverters) of the sensor can be used directly from the standard cell library of a process design kit (PDK). It makes it easy to implement the sensors directly



Fig. 4.37 Embedded sensor in a scan chain cell in support of both close-loop and open-loop sensor readouts.

with the conventional PnR tools. Figure 4.36 shows an example design which is directly placed and routed with the synopsys IC Compiler in a 28/30nm technology node. The design is relative bigger compared to the custom designed one shown in Figure 4.35, but this is expected. In this case, the area can be further reduced by decreasing the utilization percentage and putting more constraints for the tool. The sensor can be used as a design IP that is no different from other cells in the standard cell library.

As many such BTI sensor IPs are expected to be placed in a design, a methodology that is able to automatically distribute the sensor is necessary. In this thesis, we present one such candidate. The compact sensor can be embedded into a scan chain cell that is for design for test (DFT) purposes. Figure 4.37 illustrates the idea, where the newly designed scan cell has the sensor embedded. Figure 4.38 is the design methodology of integrating the



Fig. 4.38 New scancell design methodology with the BTI sensor embedded in a Synopsys design environment. .CEL and .FRAM are layout formats required by the Synopsys tools. Similar methodology can also be adapted to Cadence design environment.

sensor with the scancell from the standard cell library in a Synopsys design environment. The design compiler is used for synthesis, the output of it is the synthesized sensor netlist which can be fed into IC Compiler for PnR and generating the sensor IP layout (in .CEL and .FRAM format for Synopsys tools). The new *Scancell* netlist is edited by adding the *Scancell* component to the synthesized sensor netlist. Then the new scan cell can be placed and routed. With this methodology, a modified version of scan cell can be treated as a building block for a scan chain.

As shown in Figure 4.37, the sensor readout can be selected by a MUX logic. It can be used in two ways. In the close-loop case, run-time solutions like adaptive or accelerated active recovery can be enabled based on the sensor outputs; in the open-loop case, sensor outputs can be scanned out to users just as how regular scan chains work so that users can conduct statistical analysis or just be aware of the wearout and recovery behaviors of internal nodes, and corresponding decisions can be made off-line based on the sensor outputs. In Appendix B, we present the detailed flow and an example of instrumenting the BTI sensors in a counter design.

4.5 EM Sensing

Compared to BTI sensing, EM sensing is more challenging because of the following reasons. Firstly, EM behavior is more complex as shown in Chapter 3. Due to the stress accumulation period, the resistance keeps almost unchanged. This indicates that EM sensors need to be ON more frequently than BTI sensors so that any EM-triggered event can be recorded precisely and timely. Secondly, EM mostly occurs in power delivery network, and it is far from the back-end-of-line, which is the logic. Thus any logic-style sensor is not able to experience exactly the same conditions as the metals on top of it. Lastly, differentiating EM wearout with other wearout phenomenon such as BTI is very challenging since they show similar degradation effects to a circuit in many cases. In this Section, we mainly review some of the existing EM sensor designs and discuss how these sensors can be applied to the accelerated and active EM recovery purposes.

Since EM increases IR-drop and can lead to load performance degradations. A simple way of detecting EM is to use a ring oscillator as the load and check the frequency degradation [239]. While this method shares many similarities to many existing BTI sensing techniques as detailed in Section 4.4.1, thus it is sometime very challenging to separate EM from other wearout-induced degradations. Ring oscillator can be a perfect sensor candidate in



Fig. 4.39 Illustration of Metal-line-based EM sensors. Multiple dimensions can be used to sense at different levels.

the case when one cares about the impact of wearout as a whole only but not about which wearout leads to the degradations. Metal-line-based sensor appeared being a better candidate for EM-specific sensing and detection [201, 77]. The main idea is to put a set of metal lines on chip so that they can experience the same or more EM than the power delivery networks or other EM-sensitive interconnects. Figure 4.39 gives an example of one such EM sensing structure that is modified based on previous work [201, 77]. They are basically a set of on-chip metal lines in parallel. Since there are inherent variations in the metal wires, multiple metals are required to get statistic results. A metal line that is kept unstressed is used as the reference metal for capturing any resistance change. This idea behind this is very similar to the metastable-element-based BTI sensor discussed in previous section. The dimension of these metals need to be picked based on the system lifetime requirement, which can be translated by resistance increase threshold such as 10%. Thus the current densities in the sensing metal lines have to be relatively higher than in regular loads in order to force a shorter lifetime, and also in order to compensate for the reduced number of sensing elements compared to regular loads. Intermediate sensing check points (e.g. 5% of resistance



Fig. 4.40 The EM-induced resistance change detection circuit (design is modified based on the circuit proposed in [77]).

increase) can be added by designing metal structures with different width values as shown in Figure 4.39. The current that flows through the sensing wires, although scaled up, has to be proportional with the current in the main circuit in order to have a high correlation. Finally local operating conditions such as temperatures need to be reflected in the EM sensors, so multiple copies of such sensors need to be placed accordingly.

In order to sense the EM degradations during run-time, a resistance detection circuitry is necessary. [77] proposed one such scheme. A modified version is presented in Figure 4.40. The current source on the left is mirrored from the current of the actual loads. It provides the stress to the sensing metal lines shown in Figure 4.39. A comparator is used for comparing the voltage levels across the sensing structures and the reference metal line. If the voltage across the sensing structure is higher, then it means the EM-induced resistance increase has reached the threshold, and the sensor is triggered. *Sense* signal stays high during the sensing periods and the reference metal line is unstressed. It is put as low only during the sensor readout periods. As can be seen from Figure 4.40, introduction of analog circuit such as comparators and on-chip resistors can lead to more design effort compared to the BTI sensor design. Also EM sensing circuitry will take more silicon area. The good news is that as EM is



Fig. 4.41 Illustrations of EM sensor usage in the accelerated and active recovery case.

usually a concern for power and ground network, the number of sensors required is much less than in the BTI case where each transistor is undergoing wearout. Also the sensing structures itself is only a set of metal lines that take relatively small area. As reported by [76], one such EM sensor with 10 stressed wires costs only $100 - 500\mu m^2$ in an 180nm technology. In addition, the metal line sensing structures can be distributed in a fine-granularity way, and the resistance detection circuitry can be shared by multiple copies of such sensing structures with the extra selection logic.

As suggested by the experimental results discussed in Chapter 3, EM has a long period of stress accumulation period (e.g. 1 year) when the resistance almost stays unchanged. Thus EM sensors can be in *Sense* mode most of the time and the output of the sensor can be sampled every few days or longer. In the second region (void nucleation) of EM, where the resistance starts to increase immediately, the sensor needs to be alerted more frequently to sense the converting point between the first and second region, and also the resistance change. Accelerated and active EM recovery techniques can be applied after resistance increase to a threshold that was set based on analysis. During recovery periods, EM sensors

Туре	Design Name	Leakage Power	Dynamic Power	Area	Performance
BTI Accelerated	Neg. Voltage Generator ¹	68.85 <i>nW</i>	64.47µW	$4300 \mu m^2$	> 66.7 <i>MHz</i>
Circuits	On-chip Heater ²	16.8 <i>nW</i>	75µW	$16\mu m^2$	-
EM Accelerated & Active Recovery Circuits	Multi-mode Assist Circuit ³	-	-	$58.24 \mu m^2$	Wake-up time~ 170 <i>ns</i>
BTI Sensors	RO-Based P/NBTI Sensors	19 <i>nW</i>	-	$92\mu m^2$	-
	Metastable Element -Based Sensors ⁴	40.64 <i>nW</i>	-	$54.1 \mu m^2$	Acq. Time 12ns
EM Sensors	Metal-line Based Sensors ⁵	-	-	$100-500\mu m^2$	-

Table 4.2 Summary of	PPA Metrics f	for different	Circuit	Components	(in 28nm	FD-SOI
technology unless specif	fied)					

¹ Corresponds to a negative voltage generator designed for generating -0.3V for BTI recovery. ² Corresponds to one on-chip heater (41-stage RO) that can generate temperature of more than 80°C, in the real use case, multiple copies of this heater will be distributed.

³ PPA metrics of this circuit depends on the load size and application, here we list the examples for 8 ring oscillators running in parallel.

⁴ The numbers presented in this table shows metrics for a 2-path version of the sensor.

⁵ Data is reported by [76] with 180*nm* technology node.

need to be checked frequently as well to avoid the potential EM issues by reverse current (as shown in Figure 3.8 in Chapter 3). This has been illustrated in Figure 4.41. In summary, metal-line-based EM sensors can serve as a check engine and indicator for EM wearout and recovery, they can be integrated as a circuit element IP for an accelerated self-healing system.



Fig. 4.42 Chapter 4 Highlights.

4.6 Conclusions

In this chapter, we presented a set of circuit structures for enabling and assisting the accelerated and active BTI/EM recovery. The key components are highlighted in Figure 4.42. On-chip negative voltage is designed for generating the negative voltage for activating the BTI recovery; on-chip tunable heater is a modified version of configurable ring oscillator which can be deployed for providing high temperature when necessary. A multi-mode EM/BTI recovery assist circuit is able to support both BTI and EM recovery simultaneously, and it can be designed based on the existing power-gating infrastructure on chip, thus less design effort and overhead are required. As even with the recovery solutions, there need to be wearout sensors which directly or indirectly indicate the levels of wearout so that these recovery solutions can be asserted or daseserted. This chapter detailed 3 different sensor designs - 2 for BTI and 1 for EM. These sensors are small and flexible, and can be adapted to multiple use sceneries. In the case of proactive recovery, sensors are for sanity check so that some wearout effects (e.g. EM) are not "overly-recovered"; in the reactive recovery cases, all the actuations rely on these sensor outputs. As BTI and EM wearout effects become increasingly critical, novel techniques that are able to mitigate them with lower overhead

are highly desirable. Table 4.2 summarizes the power, performance and area metrics for the circuit components discussed in this chapter. It shows that most of these components are small in terms of area and fast in terms of response time. In summary, this chapter provides a set of circuit IP blocks and infrastructures for designing future accelerated self-healing systems. With these components at the circuit level, smart design decisions at higher levels of a system stack can be made to leverage the potential overhead and enable to reliable system. In the next chapter, we will discuss some of these higher level solutions.

Part of the work presented in this chapter has been published in [J2], [C7], [C10], [C11], [C13].

4.7 Acknowledgements

This part of work was funded by NSF CCF-1255907, SRC Global Research Collaboration (GRC) Program (Task ID. 2410.001), and C-FAR, one of six SRC STARnet Centers, sponsored by MARCO and DARPA. I would like to thank Dr. Kaushik Mazumdar (University of Virginia) for the negative voltage generator part, and thank former HPLP members for the early work in metastable-element-based BTI sensors part.

Chapter 5

Accelerated Self-Healing as a Key Design Knob for Cross-Layer Resilience

5.1 Overview

In the last chapter, circuit level solutions for accelerating and activating BTI/EM recovery were discussed. As many of these single-layer solutions are not free and they need to be triggered in a smarter way, dealing with wearout issues and applying the accelerated and active recovery need to cross layers, where various techniques are necessary to be implemented - from device level up to the application level - to work together to achieve the optimal lifetime and acceptable wearout levels with a low cost. The notion of "Cross-layer Resilience" was firstly introduced to the computing community around 2009 [156], the key idea of it is to divide error and variation tolerance into a set of tasks, which can be implemented at different levels of a system stack as listed in Figure 5.1. These resilience tasks can be treated as steps that the system follows to handle a particular reliability effect even they may not occur sequentially [38, 144]. Examples of such techniques could be cross-layer error predictions or detections for soft errors [44, 144], run-time sensing and actuation [180] and so on. Orthogonal to what has existed, the periodic accelerated self-healing techniques



Fig. 5.1 Illustration of Concept of Cross-Layer Techniques.

demonstrated in previous chapters are able to "repair" the wearout completely and efficiently by fully taking advantage of the unique recovery and frequency dependence behaviors, thus it can be a promising candidate cross-layer solution for optimizing the system resilience furthermore. In this chapter, we will present a set of solutions at different levels of system hierarchy, the circuit building blocks discussed in Chapter 4 serve as key infrastructures for the cross-layer accelerated self-healing (CLASH) system which instruments the recovery from the circuit level to the system level. The system benefits from the CLASH by operating for a longer time with higher performance in a refreshed mode, thus it will eventually lead to the significant reduction of the guardbands, and better cumulative metrics (e.g. average performance) as well. Since the notion of accelerated self-healing is orthogonal to the existing adaptive solution, so there are also good opportunities to combine them together to achieve the optimal cross-layer resilience in a more effective way with much lower cost. We will show that these CLASH techniques can contribute to provide an effective solution against wearout for both system designers and circuit designers in the design process.

5.2 Accelerated and Active Recovery Space Exploration

In Chapter 2 and 3, we have demonstrated experimentally that both BTI and EM wearout can be almost fully recovered by tuning metrics such as stress voltage/current, temperature or both. While when it comes to the actual on-chip implementations, the tunability of these metrics such as voltage range and temperature range are very limited, thus we perform a design space exploration of accelerated and active recovery for different scenarios, also we discuss solutions for finding the optimal operating schedule for both BTI and EM wearout. Details will be reported in the following sections.

5.2.1 High Temperature or Negative Voltage? or Both?

As shown in Section 2.4 and Section 2.9, high temperature of 110° C and negative voltage of -0.3V are picked and combined as effective accelerated self-healing for BTI to achieve significant metric improvements (recovery percentage, design margin reductions and average frequency improvement). While this combination shouldn't be fixed since different applications and physical locations on chip will have different available high temperatures or negative voltage resources, so it is necessary to explore other possible combinations so that they will offer the flexibility for cross-layer on chip implementations (will be discussed in the following section). Figure 5.2 shows the explored accelerated self-healing space based on the proposed model in Section 2.4.2. The same model parameters used in Section 2.7 are used here. We define "High Temperature" as the temperature above 20°C and "Negative Voltage" as any voltage values below 0V. Assume that the chip is stressed under high temperature of 110°C for 24 hours, the recovery percentage of 72.3% demonstrated in Section 2.4 is



Fig. 5.2 Accelerated self-healing space exploration for a case when a 6-hour recovery follows a 24-hour accelerated stress. Accelerated Recovery - high temperature, Active Recovery - Negative voltage.

defined as the upper limit for the reversible wearout in this study. The surface plot presents the recovery percentage (*recovered delay/net delay* increase due to wearout) under different accelerated & active recovery conditions for a 6-hour recovery period (long enough for reversible wearout to be fully recovered as shown in Figure 2.23). It turns out that the same recovery percentage can be achieved with multiple combinations of high temperature and negative voltage. For example, to achieve a recovery percentage of ~30%, the combinations could be $(50^{\circ}C, -0.3V)$, $(60^{\circ}C, -0.25V)$, $(70^{\circ}C, -0.2V)$, $(80^{\circ}C, -0.15V)$, $(90^{\circ}C, -0.1V)$, $(100^{\circ}C, -0.05V)$ or $(110^{\circ}C, 0V)$. One observation is that the resulted recovery percentage by every $10^{\circ}C$ of temperature increase is almost comparable to what is achieved by reduction of -0.05V of the voltage. This offers the flexibility and possibilities of even "controlling" the recovery levels via voltages and temperatures. It also indicates the potentials of implementing the techniques under different part of the chips (with different available high temperatures



Fig. 5.3 Illustration of the "training" process for finding the optimal stress vs. recovery balance to fully recover from the BTI wearout effects.

and/or negative voltages) while achieving the same recovery percentage even when the on-chip heat is not uniformly distributed or have fluctuations. Implementation details at different layers will be discussed in the following sections.

5.2.2 Right Balance of Wearout and Recovery for BTI

Another key factor for full recovery of BTI and EM is to employ a right balance of wearout and accelerated and active recovery so that the circuit could keep being active as long as possible but can be rejuvenated back to the fresh state within a very short sleep duration. Based on experiment results shown in Chapter 2 and 3, wearout is a relatively slow process under normal operating condition (without being accelerated by high temperature and/or voltage/current), even the reversible part of the wearout usually takes days, so the sleep period could be scheduled roughly in a daily (or several days) base - the accelerated rejuvenation follows a 1-day (or several days) period of being active.

The optimal scheduling strategy can be determined based on modeling and simulations, or it can be explored during run time for better accuracy. To actually detect the optimal



Fig. 5.4 Illustration of finding the optimal stress vs. recovery balance for EM during run time.

balance of active vs. sleep for certain system, small embedded circuit-level wearout sensors discussed in Chapter 4 need be spread over the on-chip test structures , these sensors can feed the degradation information back to the system scheduler which enables the accelerated rejuvenation techniques. Shown in Figure 5.3 is a training-like process for BTI wearout, during which the accelerated self-healing techniques can be applied incrementally (e.g. 1 day, and then 2 days) during the initial lifetime. At the time when irreversible wearout starts to show up, the optimal active and accelerated rejuvenation duration could be finally determined. This process will be able to find exactly the optimal operating schedule for the reliable-critical systems, and this schedule can be integrated as part of a system-level scheduler for recovering from wearout during the rest of lifetime.

5.2.3 Right Balance of Wearout and Recovery for EM

As EM-induced resistance increase starts after a long time, finding the optimal scheduling strategy during run time for EM is slightly different from BTI. Illustrated in Figure 5.4 is one suggested solution, in the early lifetime during the stress accumulation period, the sensors can be triggered on a monthly base to check the resistance change until it starts increasing. Then stress and recovery can be scheduled on a daily base, two day base and so on. The

optimal stress and recovery schedule is decided once the irreversible component of EM starts to be captured by the EM sensors. Compared to BTI, the overhead of finding the optimal schedule for EM is much less since during the active recovery period, the circuit blocks can still function if the circuitry discussed in Chapter 4.3 is implemented. The long period of "constant resistance" also alleviates the sensor tracking overhead.

In summary, this section and last section provide several run-time solutions to find the exact operating schedule for full recovery. On one hand, this schedule is useful for balancing the overhead of implementing recovery. On another hand, the process of detecting this schedule increases the complexities of operating a system, thus we will present another scheduling solution called "proactive recovery" given that the exact optimal stress and recovery schedule is usually unknown. Details will be covered in Section 5.4.1.

5.3 Architecture-level Accelerated Self-Healing

In Chapter 4, circuit implementations for accelerated self-healing were presented, while circuit-only solutions require very fine-grained control and can introduce huge overhead in many cases, dividing the recovery tasks across the system stack is more efficient. In this section, we will present several potential architecture level accelerated healing techniques which can be employed to leverage the tradeoffs and reduce the overhead.

5.3.1 Architectural-level Model for Wearout and Lifetime Analysis

Architecture research relies on many abstracted models that are used for capturing the lower level (circuit and device level) behaviors, these models need to be "light" in terms of simulation complexity but "heavy" in terms of accuracy. In this section, we use an opensource tool for architecture reliability called "OldSpot"¹ that is able to analyze wearout and estimate the lifetime of a system at the unit level with arbitrary distribution of workloads, and therefore variation of temperature, voltage, frequency, and so on, across time and space and arbitrary tolerance for failure. This model differs from previous work by relaxing assumptions about the behavior and failure tolerance of the system and enabling fast lifetime and reliability modeling. OldSpot is compatible with other architecture level framework, it works by receiving performance data from gem5 [24], power data from McPAT [125], temperature data from HotSpot [83], and a floorplan created with ArchFP [56], and the tool computes per-unit aging rates and runs Monte Carlo simulation to output reliability distributions for each unit and the system overall. This will enable an architect to predict the impact of expected workloads on a design in the early design phase and see which regions of the system are experiencing aging "hot spots" and should be targeted for aging mitigation techniques such as accelerated self-healing techniques discussed in this thesis.

5.3.2 Unit-level Accelerated Self-Healing

Since both BTI and EM wearout depend on temperature, voltages and switching activities, this will lead to that different units experience very different wearout behaviors during run time. There are two alternatives of instrumenting the accelerated self-healing solutions discussed in Chapter 4, the first one is coarse-grained solution where wearout sensors and recovery circuitry are distributed evenly at the core level so that the whole core experiences similar recovery behaviors, the benefit of this solution is that it is easy to control and instrument, but the downsides are that this can lead to that some units which experienced more wearout don't get the full recovery and other way around. Thus the overall benefits of recovery are still very limited considering the whole system. An more effective and

¹The development effort of this tool was led by Mr. Alec Roelke, who is a PhD student at the University of Virginia. We worked together on this project, my main contributions were: 1) Providing the device-level wearout models and parameters; 2) Helping with the idea development of structure replication; 3) Helping evaluate the accuracy of the tool.

Parameter	Value		
Instruction set	x86		
Microarchitecture	Nehalem		
Technology size	65 nm		
Supply voltage	1.1 V		
Core count	4		
CPU clock frequency	2.66 GHz		
Instruction cache	32 kB		
Data cache	32 kB		
Private L2 cache	256 kB		
Shared L3 cache	8 MB		

Parameter	Value		
Instruction set	x86		
Microarchitecture	Nehalem		
Technology size	65 nm		
Supply voltage	1.1 V		

Table 5.1 Simulated System Parameters



Fig. 5.5 Heat (a) and Wearout (b) maps representing the average temperature and relative wearout rate, respectively, when running cholesky [23]. In (a), red indicates hotter temperatures and blue indicates cooler ones. In (b), red indicates faster wearout while blue indicates slower wearout.

economic solution would be unite-level accelerated self-healing, in which "wearout hotspots" or wearout-critical units are predicted, and the accelerated self-healing is only instrumented

and applied to these units. This is also called fine-grained accelerated self-healing. With the pre-RTL reliability simulator such as OldSpot, we are able to conduct such analysis. Figure 5.5 shows an example where we run a benchmark *cholesky* from the PARSEC suite [23] on Intel's Nehalem architecture with parameters shown in Table 5.1. Figure 5.5a shows the heat map and b shows the corresponding wearout map. Here wearout is mainly dominated by BTI and EM effects. This study demonstrates that wearout is not uniform across the whole system, some functional units experience more wearout due to higher activity that leads to higher power consumption and higher temperatures. As a result, recovery solutions should be instrumented differently across the chip. For instance, based on the results from Figure 5.5b, reorder buffers (rob) experience more wearout, thus recovery circuitry such as power gating and on-chip heaters should be firstly instrumented on these rob units, and for units such as L2 cache, as wearout has a relatively small impacts on them, it is less economic to accelerate the recovery for them. In summary, unit-level accelerated self-healing provides fine-grained control at the functional block level, and it is potentially more effective because the most wearout-critical units are always implemented with more recovery resource and are also recovered first. This method is also less costly in terms of hardware overhead compared to the coarse-grained solution.

5.3.3 Dark Silicon and Core Redundancy

In multicore or Network-on-Chip (NoC) systems which consist of hundreds or even thousands of identical cores, due to the TDP limitations, a significant amount of cores cannot be operated at the same time, and this leads to the so-called *dark silicon* problem [55]. Recent research [79] has pointed out that even with the latest FinFET technology and novel processor architectures, dark silicon issue still exists and stays as a big challenge. The dark silicon usually leads to some "redundant" core resources, and these resources can be a single core or a subset of the core. Existing works tried to leverage the redundancy by developing



Fig. 5.6 A potential self-healing solution in a multi-core system. Dark silicon and core redundancy can be utilized to improve the lifetime of the whole system. t1 and t2 are two time points during the lifetime.

frameworks considering the amount of work and temperature variations and analyzing the the different redundancy arrangements [82], on top of these techniques, if these resources are scheduled and allocated in such a way that they can be healed by the generated heat from the active elements, the average lifetime of the whole system will be significantly improved.

It has been shown in Figure 5.5a that the temperature difference between the active regions and the inactive regions can be as large as 30°C, similar observations were also made in alpha processor and other modern many-core designs [241]. Thus the huge amount of generated heat by the active elements can be fully took advantages of. Figure 5.6 illustrates a potential implementation in a simplified multicore system, where the sleeping cores (or resources) are always surrounded by the active cores (or resources), by switching the workload among different cores, these sleeping cores can be deeply rejuvenated by the heat generated by the

active cores. This method is application dependent and also needs the support of the system scheduler. There will also be switching overhead, but this solution doesn't need to stall the system completely during recovery, and it also overcomes the power overhead introduced by circuit level on chip heaters (presented in Chapter 4). Thus this method could be potentially used together with the circuit level solutions and system scheduling solutions to leverage the tradeoffs. Details of scheduling will be discussed in Section 5.4, and a potential cross-layer implementation will be covered in Section 5.6.

5.4 Scheduling at the System Level

Right balance of wearout and recovery can lead to almost full recovery for both BTI and EM recovery, and this has inspired the idea of active recovery at the system level, at which task scheduler can implement scheduling based on the desired "circadian rhythms", this section will briefly discuss several ways of doing scheduling and how to apply these in different applications.

5.4.1 Reactive Recovery vs. Proactive Recovery vs. AC Stress

Since recovery can be used as an effective to deeply rejuvenate the system from wearout, it will be worthwhile to schedule the recovery periods proactively and during these recovery periods several accelerated and active recovery techniques can be applied. But it is important to decide when to insert these recovery due to the irreversible component of wearout. Philosophically, there are two alternatives for scheduling recovery: reactive, when a system or part of system has aged by a particular threshold amount, and proactive, in anticipation of future wearout. Figure 5.7 illustrates the differences. In the reactive recovery case, a wearout threshold a% is preset as the upper limit and starting point for the recovery periods. This threshold can be pre-calculated based on the wearout models or simply can be the design



Fig. 5.7 Illustration of Reactive Recovery vs. Proactive Recovery vs. AC Stress.

margin. During operation, wearout sensors need to be able to track and feed the wearout information back to the scheduler frequently, as it should be noticed that the time length to reach the threshold is unpredictable due to the different switching behaviors during different periods of lifetime. In this case, wearout sensors are very important and critical for deciding when to start recovery. As for the proactive recovery case, the recovery schedule is preset, and this schedule needs to be set such that the irreversible components haven't started to accumulate. Thus it is determined based on pre-design space exploration (e.g. with OldSpot) and wearout models which are usually provided by the foundry. Once this schedule is set, the scheduler inserts the recovery periods proactively so that recovery can always compensate the accumulated wearout effects. In this case, wearout sensors are still necessary, but just for sanity check purposes. Proactive recovery acts similar to a "slow AC stress" with a skewed

duty cycle, in which recovery and stress appear as pairs. Note that it is not exactly identical to normal AC stress (defined in Chapter 2.6.1) which is also shown in the figure, AC stress is when there are switching activities and transistors/metal wires are under periodic stress and passive recovery only, wearout still accumulates in these cases as passive recovery has been demonstrated being slow.

Reactive accelerated recovery seems more "economic" since it is only scheduled when needed. But it needs to track changing wearout levels (such as threshold voltage increase or resistance increase), has the disadvantage of being unpredictable, thus potentially introducing performance and/or energy overheads at inopportune times and likely leading to a smaller improvement in lifetime, and accumulates upfront more irreversible aging thus leading to a lower expected performance and energy – circuit operates more time in wearout mode. In addition, the preset wearout threshold is hard to determine. While for proactive recovery, the recovery periods are scheduled ahead of any sign of wearout, is easier to implement and results in the system operating for longer time in a "refreshed" mode, thus leading to better expected performance and other metrics, and has better cumulative metrics as well. Most importantly, the preset schedule can even follow the users' schedules and be fully compatible with human "circadian rhythm". Details will be discussed in the following section.

5.4.2 Application-dependent Proactive Recovery with Scheduling

Based on experimental results and analysis in Chapter 2.9.2 and 3.6, wearout is a relatively slow process under regular operating conditions (without being accelerated), even the reversible part of the wearout usually takes at least one day under the constant stress (e.g. 31*hours*), so the sleep period could be scheduled roughly in a daily (or several days) base – the accelerated and active recovery follows a 1-day (or several days) period of being active. Here we define two use cases and discuss the possible scheduled patterns correspondingly.



Fig. 5.8 Recovery time under different accelerated & active recovery conditions after 12-hour constant stress under regular operation condition (room temperature, nominal V_{dd}).

• *Use case 1*: For applications like mobile devices and consumer electronics, the scheduled active vs. sleep pattern could even follow human beings' circadian rhythms, for example, the devices are active during daytime for 12 hours, then during sleep, although certain blocks need to be still active for data retention, some critical performance-hungry components could be deeply rejuvenated for next day's "full speed". The length of the recovery time depends on the conditions (temperature and voltage). Assume that the system or is active for 12 hours under regular conditions (room temperature and nominal voltage). In the worst case, it is always under constant (DC) stress, but still in irreversible wearout region (within 31 hours). Based on the model in Chapter 2.4.2

and analysis in Chapter 2.9.2, the required time for full recovery under different sleep condition is plotted in Figure 5.8. It shows that recovery is very fast and needs about 23.2 minutes under the best combinations of 110° C and -0.3V. For a slight higher temperature or negative voltage, the recovery can be much faster than the regular operating condition. For example, under 50°C and -0.3V, it only needs about 57 minutes to fully recover to the fresh state. These devices could benefit from such short sleep intervals through periodic rejuvenation.

• Use case 2: For data center or server applications, although the system runs most of the time, and it might not be feasible to fully turn off the system as frequently as in mobile devices or embedded systems mainly because of the high setup times and full usage of the system. While recent research [27, 64, 160] have proposed energy-efficient solutions that implement novel load balancing and/or scheduling algorithms so that the idle and lightly-loaded cores are able to be switched to sleep states. For example, in [64], authors introduced a dynamic power management policy called *SoftReactive*, which is able to put servers to sleep in a conservative way by setting a wait time when load drops, and when the load increases, it turns servers back. This work demonstrated that as the scale of data center increases, the effectiveness of sleep states is even more pronounced because the setup time has less and less effect on performance. During these intrinsic sleep state, circuit level and architecture level accelerated self-healing techniques (discussed in the last section) can work together to push for the full recovery capabilities. Larger data centers can exploit more benefits of sleep states for recovery even by employing a simple dynamic power management policy.



Fig. 5.9 An illustration of recovery-driven design methodology in an IC design cycle.

5.5 Recovery-Driven Design Methodology for Resilient System Design

It has been discussed in previous sections that accelerated self-healing can be achieved at different system hierarchy during run time, but the hardware blocks and all the scheduling decisions need to be instrumented during design phase. For example, where to place the sensors, and how many of the recovery circuitry are required, and what will be the tradeoffs. In this section, we will discuss a potential design methodology that is driven by recovery.

Figure 5.9 illustrates the recovery-driven design flow. It starts from the design specifications (SPEC) and applications which define the expected lifetime, clock frequency, power and area budget and so on. This flow assumes that high level architecture level decisions based on the SPEC have been made, these decisions can be type of ISA, number of cores and so on. These information can be used as the inputs for pre-RTL and architecture level simulation tools to estimate the power, thermal and performance behaviors. Reliability tools such as OldSpot discussed in Section 5.3.1 that are built based on the wearout models can take the reported behaviors from other tools and estimate the wearout behaviors at a higher level, such as which units could be potentially wearout-critical and what would be the possible operating schedules for the proactive recovery. These decisions can guide the RTL design phase when designing the microarchitectures, for example, more wearout sensors can be instrumented for those wearout-critical units, multiple copies of these units can be integrated in the design, and so on. The control logic for closing the wearout and recovery loop is also implemented in this stage of the design process. By instrumenting the recovery during logic design, it will build the infrastructures for the run-time recovery solutions as discussed in previous sections. During the physical design stage, core/resource allocation solutions that are beneficial to accelerated recovery (discussed in Section 5.3.3) and PDN design against wearout (discussed in Chapter 4.3) can be part of the flooplan decisions. Similarly, the recovery circuitry such as on-chip heating elements or sensors can be placed close to the wearout-critical units. During the rest of the physical design steps, circuit knobs that affect the recovery levels will also be considered during the physical design, examples of such knobs are logic depth (discussed in Chapter 4.2.2) or power gating styles (discussed in Chapter 4.2.3). As reliability is usually not addressed in every design stage in most of the current design flows, this section gives an


5.6 Putting It All Together – CLASH: Cross-layer Accelerated Self-Healing System 155

Fig. 5.10 An illustration of Cross-Layer Accelerated Self-Healing (CLASH) System.

example of considering reliability during the whole design process, and it contributes to the efforts on developing future design methodologies for ensuring reliability.

5.6 Putting It All Together – CLASH: Cross-layer Accelerated Self-Healing System

In previous sections, circuit, architecture and system-level accelerated self-healing solutions were discussed. The recent shift of architecture to heterogeneous and many-core systems significantly increases the number of integrated cores. Specialized computing resources serve for different load tasks, which also leads to different EM and BTI behaviors, thus requiring different recovery strategies. This requires the designers to integrate single layer solutions all together in a cross-layer way to minimize the cost for ensuring reliability.



Fig. 5.11 Illustration of periodic proactively scheduled EM/BTI recovery.

Figure 5.10 presents one potential cross-layer accelerated self-healing (CLASH) system with the localized active recovery techniques. As discussed in Section 5.3.2, at the architecture level, localized active recovery at the core level or unit level will be able to leverage the cost while rejuvenating the "aged" system. Each red/yellow square represents a core or logic block with local PDN and can have different recovery strategies. Blue squares refer to the recovery circuitry and wearout sensors that are distributed across the unit. In the meanwhile, "Dark Silicon" still appears at a big challenge in these systems. The "dark" parts of the chip usually lead to some "redundant" resources which have intrinsic OFF periods, and these resources can be a single core or a subset of the cores. Since we have demonstrated that high temperature is able to accelerated the recovery of both wearout mechanisms, and if these redundant resources (e.g. the core located in the center of the system shown in the figure) can be scheduled and allocated in such a way that they can be healed by the generated heat from the neighboring active elements, the recovery can be further sped up (discussed in Section 5.3.3). Figure 5.11 presents an example of run-time proactive scheduling for BTI and EM active recovery. In the early lifetime, since EM-induced stress hasn't reached the nucleation threshold, the performance degradation will be caused mainly by BTI; novel BTI and EM sensors can be employed to track wearout and check the degradation information. Short

intervals of BTI active recovery periods are inserted proactively to bring the chip back to the fresh status in time; during these intervals, certain states need to be in retention mode, alternatively, workload can be shifted to other redundant resources. EM active recovery period can be scheduled either from when the void nucleation happens or even earlier. Based on the measurement results presented in Chapter 3, early EM recovery is more economic and efficient, and the system is still in operation during EM recovery interval, so EM active period can be scheduled alternately during normal operation with a small switching overhead. Overall, such a scheduling strategy can potentially fully recover both the BTI and EM wearout, such that the system always runs in a "refreshing" mode; the necessary wearout guardbands can then be significantly reduced as well.

5.7 Tradeoff Analysis

Accelerated self-healing techniques are not free in terms of hardware cost, but we believe that through the cross-layer implementations, the tradeoff among power, performance and reliability can be balanced. Some of these tradeoffs can be leveraged during design time as described in Section 5.5, and some need to be integrated into the run-time solutions. Shown in Figure 5.12 is a suggested implementation flow. For example, generating heat only with on-chip heating elements or providing the negative voltages by the circuit-level voltage generator together will cause additional power overhead during the sleep intervals, and this can be expensive especially for those energy-constrained applications. To avoid this cost, at the architecture level, the core allocator or load balancer can first work together to assist the process by shifting the applications among different cores and making the right allocations (e.g. as shown in Figure 5.6) so that the circuit level solution can "rest". While due to the existing of the cooling system, the generated heat might be released soon, thus the circuit level heating and negative voltage generator could restart again. At the upper hierarchy, the system scheduler leverages the tradeoffs and schedules the necessary accelerated self-healing



Fig. 5.12 Cross-layer implementation of accelerated self-healing for leveraging the tradeoffs. At the circuit level, recovery circuit and wearout sensors are distributed on the wearoutcritical units, and they are triggered when needed. Architecture level accelerated self-healing solutions can utilize some intrinsic sleep behaviors and heat to recover inactive parts. It can also compensate some of the power overhead introduced by the circuit level recovery solutions. System level scheduling is able to divide the recovery tasks and make the high-level recovery decisions.

periods proactively (e.g. following a pattern as discussed in Section 5.4.1). Key elements that enable the transparency between each layer are wearout sensors. Such sensors can be physical sensors (BTI sensors in Chapter 4.4 and EM sensors in Chapter 4.5) at the circuit level, program counters at the architecture level, or virtual sensors at the system or application levels. This implementation flow requires hardware and software to work together in a similar way to other cross-layer resilience solution, but the collaborative efforts from all layers working together can potentially ensure that the whole system will always be rejuvenated and run reliably over the entire lifetime with a relatively low cost.



Fig. 5.13 Chapter 5 Highlights.

5.8 Conclusions

As implementing accelerated self-healing at a single layer can potentially lead to huge power and area overhead, we discussed the cross-layer implementations of the accelerated self-healing in this chapter. As highlighted in Figure 5.13, we discussed both architecture and system-level recovery solutions that can benefit from the accelerated self-healing. Distributing the recovery tasks across the system stack can improve performance, power and area costs by taking advantages of the strengths of each layer. We also commented at a high level how recovery can be integrated as part of the design decisions in a conventional design cycle. Overall, accelerated self-healing, as a promising technique for wearout behaviors, can be potentially introduced as a key design knob for cross-layer resilience during the design process to achieve the optimal resilience level effectively.

The work discussed in this chapter has been published and presented in [J2], [C2], [T4] and [P4].

5.9 Acknowledgements

This part of work was funded by NSF CCF-1255907, SRC Global Research Collaboration (GRC) Program (Task ID. 2410.001), and C-FAR, one of six SRC STARnet Centers, sponsored by MARCO and DARPA. I would like to thank my colleague Alec Roelke (from HPLP lab, University of Virginia) for developing the architecture level reliability simulator "OldSpot" and for generating the heat map and wearout map for Chapter 5.3.2.

Chapter 6

Exploring Circuit Aging in FinFET-enabled Internet of Things (IoT) Applications

6.1 Overview

Over the past few decades, there have been continuous speculations for demise of Moore's law [145] due to the physical limit of transistors and manufacturing complexity. FinFET technology has appeared as a result of the relentless increase in the levels of integration, and they allowed scaling below 20*nm*, thus helping to extend Moore's law by a precious decade with another decade likely in the future when scaling to 5nm and below. Due to superior electrical parameters and unique structure, these 3-D FinFET transistors offer significant performance improvements and power reduction compared to planar CMOS devices. As we are entering into the sub-10nm era, FinFETs have become dominant in most of the high-end products, and it has been also an ideal candidate for low power and energy-constrained applications, of which Internet of Things (IoT) kept gaining attentions in the last decade. IoT brings a paradigm where humans and "things" are connected and has been a powerful enabler

to make technology more human centric and real time. Although the IoT industry has not fully migrated to deeply scaled technologies because of cost and leakage issues but recent advances in technology such as FinFET provide a compelling combination of performance, power, highest integration and ease of design for low-power IoT products. This has pushed the IoT industry to start adapting to these advanced technology nodes [133].

As the transition from planar to FinFET technologies is still ongoing, it is important for digital circuit designers to understand the challenges and opportunities brought in by the new technology characteristics. As the first part of this chapter, we will study these aspects from the device to the circuit level, and we make detailed comparisons across multiple technology nodes ranging from conventional bulk to advanced planar technology nodes such as Fully Depleted Silicon-on-Insulator (FDSOI), to FinFETs. In the simulations we used both state-of-art industry-standard models for current nodes, and also predictive models for future nodes. Our study shows that besides the performance and power benefits, FinFET devices show significant reduction of short-channel effects and extremely low leakage, and many of the electrical characteristics are close to ideal as in old long-channel technology nodes; FinFETs seem to have put scaling back on track. However, the combination of the new device structures, double/multi-patterning, many more complex rules, and unique thermal/reliability behaviors are creating new technical challenges. Adapting to the new challenges and fully benefiting from FinFETs will require the growing knowledge and design experiences and this part of the work contributes to add to that knowledge base.

From application perspectives, reliability of the IoT devices becomes extremely critical. Circuit aging ¹ discussed in this thesis has a direct impact on lifetime of these devices and their performances. As IoT becomes a general-purpose technology which starts to adapt to the advanced process nodes, it is necessary to understand how and on what level aging affects different categories of future IoT applications. In previous chapters, we mainly investigated

¹This chapter focuses on circuit aging (especially BTI and HCI) only. Battery aging and socket (and holder solder) aging in IoT devices are also important, but they are out of the scope of discussion.

the mechanisms of circuit aging issues and how they can be recovered. In the second part of this chapter, we will look into impact of the circuit aging in the context of applications, especially in IoT domain. We will answer these questions by conducting extensive circuitlevel simulations with foundry-calibrated transistor aging models in advanced FinFET node. Since aging is highly dependent on application behaviors that define the operating voltage, temperature and active time, we perform a survey of existing IoT applications and classify them based on aging-related metrics. By studying aging behaviors in each category, we show that aging can impose immense degradations in performance and design margin for some IoT applications. Our results prove that to meet IoT lifetime requirement, both battery lifetime (energy-efficiency perspective) and chip lifetime (circuit aging perspective) need to be considered together in the full design cycle. As flat guard-band approach could introduce unacceptable energy overheads for IoT systems, several dynamic solutions that are able to take advantage the recovery behaviors (covered in previous chapters) are also presented in this chapter.

6.2 Back to the Future: Digital Circuit Design in the Fin-FET Era

6.2.1 Motivation for Studying FinFET Technology

The continuous scaling of planar CMOS devices has delivered increasing performance and transistor densities. However, it also reached a point where increased leakage current, fluctuation of device characteristics and short channel effects became serious obstacles to further scaling. This was mainly because deeply-scaled planar devices became increasingly influenced by the drain potential as the gate lost the ability to fully control the channel; this led to transistors that were never fully off and leaked continuously. To solve this problem, gate oxides were aggressively thinned and high-k dielectric gate materials were adopted to

increase the gate-channel capacitance, but the gate-related issues, such as gate leakage and gate-induced drain leakage (GIDL) increased [21, 33]. FinFET devices became attractive for sub-30nm nodes [226, 7] because of their unique channel structure with good gate control that enables a much improved short channel control, thus requiring little or no doping in the channel. The threshold voltage V_{th} can be scaled down in FinFETs for both improved device performance and a much lower operation voltage. Lower channel doping also reduces dopant ion scattering, thus leading to better drive currents and decreases random dopant fluctuations (RDF) [117, 236, 101]. FinFETs back-end-of-line (BEOL) fabrication is fully compatible with planar devices in both bulk and SOI varieties, which reduces the need for new, FinFET-specific developments in that area. However, the introduction of FinFETs has brought a few changes and challenges in digital circuit design due to their unique gate structure and electrical properties. This has also impacted the circuit design decisions and some of the available design tradeoffs. For example, FinFET devices have a significant amount of parasitics that need to be modeled precisely and be carefully considered in the layout of all circuits, especially in SRAM and analog circuits. From a circuit design aspect, in addition to the extra effort needed to address the impact of parasitics at the layout level, new circuit techniques are needed in the area of body-biasing and memory read/write assist in SRAMs to replace techniques that worked well in planar but are inefficient for FinFET. The double/multipatterning also requires tool vendors and designers to work together to make sure the layout coloring is correct (colors refer to different exposures of the same layer while performing multipatterning). New constraints have been added to FinFET design, such as width quantization and self-heating effects, for which designers need to make early decisions in the design cycle. In this chapter, we analyze these aspects at both the device and circuit levels. To study these challenges, we simulate across multiple technology nodes which cover a wide range of gate lengths and also substrates including both FD-SOI and



Fig. 6.1 Illustration of Structural Differences (No substrate): (a) Planar Device; (b) FinFET Device.

Bulk. For FinFET, we simulate with both 1xnm industry-standard node² and a 7nm predictive node [48]. We restrict our focus to digital circuits, but several of the findings can be applied to analog design as well.

6.2.2 FinFET Scaling and Sizing

Compared to conventional planar devices (bulk or SOI), FinFET devices have unique 3-D gate structures that enable some special properties for FinFET circuit design which will be detailed in the following sections. Illustrated in Figure 6.1 is a planar device and a FinFET device (the substrate is not included in the figure). While the channel of the planar device is horizontal, the FinFET channel is a thin vertical fin with the gate fully "wrapped" around the channel formed between the source and the drain. The current flows parallel to the die plane whereas the conducting channel is formed around the fin edges. With this structure,

²In advanced technology nodes the "numbering" scheme is somewhat arbitrary, while in older technologies the node "number" used to denote the smallest feature size, usually the transistor gate length, in modern technologies the node number does not refer to any one feature in the process, and foundries use slightly different conventions; In this chapter, we use 1x to denote the 14nm-16nm FinFET nodes offered by several foundries.

the gate is able to fully deplete the channel thus having much better electrostatic control over the channel.



Fig. 6.2 Cross section view of structural differences between (a) Bulk FinFET and (b) SOI FinFET.

FinFETs can be classified by gate structure or type of substrate. Different gate structures lead to two versions of FinFET - Shorted-gate (SG) FinFETs and Independent-gate (IG) FinFETs. In SG devices, the left and right sides are connected together in a wrap-around structure as in Figure 6.1; this can serve as a direct replacement for the planar devices which also have one gate, a source and a drain (three terminal- devices). In IG FinFETs, the top part of the wraparound gate structure is etched out and this results two separate left and right sides that can act as independent gates and can be controlled separately [142, 148]. Although IG FinFETs offer more design options, the fabrication costs are also higher in general. Depending on the substrate, the FinFETs can be either SOI or bulk FinFETs as illustrated in Figure 6.2. SOI FinFETs are built on SOI wafers and have a lower parasitic capacitance and slightly less leakage. Bulk FinFETs are more familiar to designers, the fabrication costs are relatively lower, and they also have better heat transfer rate to the substrate compared to SOI FinFETs [142], thus bulk FinFETs are usually preferred for most digital applications. The fabrication of both types of FinFET devices is compatible with those of the conventional planar devices fabricated on either bulk or SOI wafers.

Unlike planar technologies for which the transistor width is a continuous value fully under the control of the circuit designer, in FinFET technologies device widths are quantized into units of whole fins. The effective gate width of a FinFET device is roughly $n(2H_{fin}+t)$, where n is the number of fins, t is the fin width and H_{fin} is the fin height as illustrated in Figure 6.1b. Since the gate of a FinFET device is designed to achieve good electrostatic control over the channel, and because of the etching uniformity requirements, the fin dimensions (e.g. height H_{fin}) are not under designer control, and thus the device width cannot have an arbitrary value as in planar technologies. Wider transistors with higher on-currents are obtained by using multiple fins, but the range of choices is limited to integer values. This is known as the *width quantization* issue [190, 71, 235]. This quantization issue doesn't allow flexibility in terms of device sizing which becomes problematic especially in analog design and SRAMs. The designers need to adapt to this new constraint during the design phase [206]. An alternative solution would be for the foundry to provide the designers with multiple versions of FinFET with different fin heights [120]. For example, [127] did an early attempt by exploring the design space of FinFETs with double fin heights and showed that the lack of continuous sizing can be somewhat compensated; this method though has many uncertainties from both fabrication costs and manufacturing difficulties, so it is unlikely to become widely available. In summary, for digital circuits, width quantization might not be a big issue since most of the cell designs can be adapted to use the limited choice of device widths available.

As fin height H_{fin} determines the overall width of a device. This is a very important parameter for circuit designers but they don't actually have control over it. Smaller fin heights offer more flexibility in terms of sizing, but this would lead to more fins, which means more silicon area. In contrast, FinFET devices with taller fins offer less flexibility with sizing but have a smaller silicon footprint and the increasing fin heights for successive FinFET nodes combines with the lateral scaling to actually accelerate "Moore's Law"-style scaling; but this

Technology	Physical Length L_g (nm)	Nominal V _{dd} (V)	I_n/I_p (Saturation)	Subthreshold Slope (mV/dec)	DIBL Parameter	GIDL Slope (mV/dec)	Channel Length Modulation λ (/V)
130nm Bulk	120	1.2	4.24	92.07	0.53	3346	0.246
45nm Bulk [63] ¹	45	1.0	1.45	98.3	1.61	286	0.387
28nm FDSOI	30	1.0	3.21	84.2	0.993	198.42	0.260
1xnm Bulk FinFET	14	0.8	0.99	71.1	0.485	429.79	0.256
7nm Bulk FinFET [48] ¹	7nm Bulk inFET [48] ¹ 20 0.7		0.90	67.6	0.745	2220.6	0.203

Table 6.1 Summary of device parameters across multiple technology nodes (extracted from I-V curves)

¹ Predictive Nodes

might also result in larger short-channel effects and some structural instabilities [21, 142]. In addition, taller devices could also lead to an increase in unwanted capacitance. This indicates that there are some opportunities for device-circuit codesign that are unlikely to become available for fabless companies but could become important for vertically-integrated companies that have their own fabs. An example of such involvement can be to analyze the design space of current versus capacitance for different fin heights. As the technology node approaches the sub-10nm scale, this type of analysis is more and more important since the fabrication difficulties are increasing, and the design tradeoffs might drastically change [59].

6.2.3 A Comprehensive Study of Bulk vs. FDSOI vs. FinFET Devices

From a digital circuit designer's perspective, whether the technology is planar or FinFET, whether it is bulk or SOI, the parameters of interest are the same - how much current can one transistor drive, leakage, DIBL, GIDL and so on. Summarized in Table 6.1 are device parameters we extracted based on extensive simulation results across multiple technology nodes. As for the 7nm node, we use a recently released predictive 7nm PDK [48] which is based on current realistic assumptions for the 7nm technology node but is not tied to or



Fig. 6.3 Capacitance components for a FinFET device: (a) Cross-section view and (b) Top view.

verified by a specific foundry. We believe that this analysis will provide us with a good insight on how FinFET devices are right now (with industry PDKs) and how good these devices are likely to be in the future (with predictive PDKs) as we move forward compared to the planar devices. In the follow sections, details of each parameter will be discussed.

Device Models

Device models are critical for circuit designers to run simulations and make design decisions. They need to be accurate and efficient in terms of simulation time and complexity. The fact that fins are 3D structures that rise above the substrate means that they are more strongly affected by their immediate environment than planar devices. This results in a number of challenges during the modeling process. For example, the interaction between the device and its surroundings needs to be accurately modeled. Besides, the unique gate structure leads to increased gate capacitance and also to more components when modeling the parasitic capacitance and resistance compared to the planar devices [175, 131]. These capacitance and resistance values are crucial since the inaccuracy caused during extracting R and C parasitic will lead to mis-characterization and under/over-estimated design margins.

Figure 6.3 shows an example of how FinFET parasitic capacitance is accounted for a 2-finger device. It is clear that more components contribute to both intrinsic capacitance (in the SPICE models) and extraction capacitance (accounted during extraction). For example, the gate capacitance includes gate to top of fin diffusion, gate to substrate between fins, gate to diffusion inside channel, gate to diffusion between fins, gate to contact, and so on. Similarly, the Fin-to-Fin capacitance is also newly introduced for FinFET devices. The complexity of modeling has been increasing as the device dimensions shrink. Coupling and Miller effects are more pronounced in these devices as well.

The FinFET structure brings new modeling challenges. In a planar device, the source and drain are self-aligned with the gate and often intrude slightly under it. In FinFET devices there is a spacer between the gate and the source and drain, which are usually raised and have a strain caused by a SiGe layer that creates a lattice mismatch. This means there are much more complex parasitic capacitance and resistance structures and more model calibrations are required to achieve good accuracy. As for the designer, the simulation efficiency also matters and it depends on the levels of model complexity, but thanks to the fast solvers and accurate extraction tools recently developed, the simulation time has remained tractable.

Leakage

One of the driving forces that leads the industry to move from bulk planar to FD-SOI or FinFET technologies is the difference in leakage. With every new process generation the doubling of gate density is also associated with a doubling of the amount of leakage current [81]. This is also clear from the simulation results in Figure 6.4 where the subthreshold current (OFF current) per unit width is plotted for different technology nodes. It can be seen from the plot that, when scaling from 130nm to 45nm, the leakage current increases significantly, due to the fact that the channel depth underneath the gate becomes larger and a significant volume of the channel is too far away from the gate and there is a subsequent loss



Fig. 6.4 Leakage current evolution with technology scaling.

of electrostatic control. FDSOI and FinFET on the other hand achieve much better leakage results because the gate has much better control over the channel in these technologies. Our simulations show that 28*nm* FDSOI and 7*nm* FinFETs have comparable leakage numbers. However, 1*xnm* bulk FinFET shows a reduction of leakage of at least 50%. This can be due to the fact that FDSOI and FinFET use different mechanisms to reduce leakage. In FDSOI, leakage reduction is achieved by making the channel thinner, by limiting its depth with the help of an insulating layer, while in FinFET it is achieved by making the gate wrap around the channel.

Another way of explaining the leakage reduction in FinFET devices is to look into the subthreshold slope. The sub-threshold slope also measures how fast the device can switch from OFF to ON, and the lower bound is 60mV/dec at room temperature. Table 6.1 shows that, together with the move to FDSOI and FinFET, the subthreshold slope value has actually improved with scaling and this has resulted in a significant benefit for continuously



Fig. 6.5 $I_{\rm on}/I_{\rm off}$ ratio with technology scaling.

improving frequency, active power, leakage power or a combination of the three over the past few years [178].

*I*on/*I*off Ratio

The $I_{\rm on}/I_{\rm off}$ ratio is an important figure of merit for having high performance (higher $I_{\rm on}$) and low leakage power (lower $I_{\rm off}$) for the devices. Since the leakage current ($I_{\rm off}$) has been significantly reduced in FinFET devices, their $I_{\rm on}/I_{\rm off}$ ratio is superior to bulk, as shown in Figure 6.5. This has also enabled a continuous performance improvement.

Drain-Induced-Barrier Lowering (DIBL)

Drain-Induced-Barrier Lowering (DIBL) is a short-channel effect that appears as the distance between the source and drain decreases to the extent that they become electrostatically coupled. The drain bias affects the potential barrier to carrier flow at the source junction, resulting in subthreshold current increase. To characterize it, we use the DIBL parameter, which is defined in Equation 6.1 and corresponds to the change of leakage current due to V_{ds} . The smaller this parameter, the better the DIBL behavior is. It is shown in Table 6.1 that FinFETs achieve very good DIBL behaviors compared to bulk devices. In particular, the 1xnm FinFET device has the lowest DIBL effect among all five technology nodes considered.

$$\Delta log(I_{\rm off}) = (DIBL \ Parameter) \times V_{ds} \tag{6.1}$$

Channel Length Modulation (CLM)

Channel length modulation (CLM) is another short-channel effect that is caused by large drain biases. It is characterized by the CLM parameter λ which is generally proportional to the inverse of the channel length. Smaller λ means less CLM effect. Table 6.1 shows that CLM has been getting worse as the channel length shrinks in planar devices even by increasing the doping density. When technology switched from planar to FDSOI and FinFET, CLM has been improved due to the better control over the channel. Especially, in 7nm technology node, the CLM effect is the smallest and is as good as a relatively old long-channel technology (130nm).

GIDL

The introduction of high-k/metal-gate stacks in planar devices has led to substantial reduction in the gate leakage and has exposed other leakage mechanisms such as gate-induced drain leakage (GIDL) as primary gate-related leakage mechanisms [104]. GIDL occurs due to the high reverse bias between the silicon body and the drain junction (a PN-junction) near the gate edge at a nearzero or a negative gate bias [47]. GIDL usually increases as the gate length (L_g) decreases due to the floating body effect and is usually pronounced in short-channel devices. In this work, we pick the GIDL slope to quantify this effect; the



Fig. 6.6 VTC curves under different supply voltages for a 1xnm FinFET inverter (PMOS and NMOS are sized equally).

larger this slope the lesser GIDL effect the device has. Interestingly, the results in Table 6.1 indicate that as the technology switched to FinFET, GIDL has actually improved. The suppression of GIDL can be explained by the light doping of the channel and better junction placement gradient as suggested in [104]. In conclusion, FinFETs are superior to planar devices in terms of $I_{\rm on}/I_{\rm off}$, DIBL, CLM, GIDL, and thus appear to be a true "back to the future" reset of most of the metrics that were getting worse with every new technology node for bulk planar technologies.

$W_{\rm p}/W_{\rm n}$ Ratio

Another interesting aspect for FinFET technologies is that the pull up network (PUN) and the pull down network (PDN³) can become very symmetric. PMOS and NMOS devices with the same number of fins have very comparable driving strength, and the conventional

³Here PDN refers to pull down network, in the rest of thesis, PDN means power delivery network.

2:1 or 3:1 sizing strategy is not be applicable (or necessary) in the FinFET case. This can be seen from the I_n/I_p ratio in Table 6.1, which is very close to 1 for the FinFET nodes. Figure 6.6 further demonstrates this. It plots the voltage transfer curve (VTC) under different supply voltages for a FinFET inverter with $W_p/W_n = 1$. It shows that the small-signal gain (which is the slope of the transfer curve when the input is equal to the mid-point voltage) is close to ideal (very high gain), and the curves are very balanced in all cases which further demonstrates that the ratio of 1:1 is optimal for FinFET logic.

The reason behind this fact is due to the unique fabrication process for FinFET. As opposed to planar structures which can only be fabricated in a single plane due to process variation and interfaces traps, FinFETs can be fabricated with their channel along different directions in a single die. This results in enhanced hole mobility. The N type FinFETs implemented along plane $\langle 100 \rangle$ and the P type FinFETs fabricated along plane $\langle 110 \rangle$ lead to faster logic gates since it combats the inherent mobility difference between electrons and holes [21, 45, 137]. Moreover, since the gate has very good control over the channel, doping concentrations can be much lower than in planar devices, thus allowing to reduce the random dopant fluctuations (RDF) [101], mitigating the impact of mobility on current. The symmetric PUN and PDN introduce ease in terms of physical design and sizing but it also brings slight changes in design decisions and standard cell design.

Alpha-power Law

The long-channel MOSFET model (Shockley model), assumes that carrier mobility is independent of the applied fields, since the lateral or vertical electric fields were low [229]. However, for short-channel MOSFETs, the velocity of carriers reaches a maximum saturation speed due to carriers scattering off the silicon lattice. This also leads to a degradation in mobility that depends on the gate to source voltage V_{gs} .



Fig. 6.7 Velocity saturation index (α) for different technologies.

The drain current I_d is quadratically dependent on the drain to source voltage (V_{ds}^2) in the long-channel regime and linearly dependent on V_{ds} when fully velocity saturated due to an electric field higher than a critical electric field $E_c = V_c/L_g$ [179], where V_c is the corresponding critical voltage and L_g is the gate length. A moderate supply voltage is when the transistor operates between the long-channel regime and velocity saturation. The complete model, called the α -power law model, is presented in Equation 6.2:

$$I_{ds} = \begin{cases} 0, & V_{gs} < V_t \quad \text{(Cutoff)} \\ I_{dsat} \frac{V_{ds}}{V_{dsat}}, & V_{ds} < V_{dsat} \quad \text{(Linear)} \\ I_{dsat}, & V_{ds} > V_{dsat} \quad \text{(Saturation)} \end{cases}$$
(6.2)

where $I_{dsat} = P_c \frac{\beta}{2} (V_{gs} - V_t)^{\alpha}$ and $V_{dsat} = P_v (V_{gs} - V_t)^{\alpha/2}$. The exponent α is called the velocity saturation index, and ranges from 1 for fully velocity saturated transistors to 2 for transistors with long channel or low supply voltage. We performed $I_{ds} - V_{gs}$ simulations for the base NMOS transistors of four different technologies and determined their respective

velocity saturation index α . The results obtained, summarized in Figure 6.7, suggest that, as we switch to FinFETs, devices behave increasingly more according to the long-channel model, again, in a "back to the future" way.

6.2.4 FinFET Fabrication

The process technology of FinFET is relatively straightforward and compatible with conventional planar device fabrication process [240]. But there are still challenges, for example, fin shape control and recess of shallow trench isolation (STI) oxide are still critical in the integration of FinFETs. This section lists several of such fabrication advances and challenges in the FinFET era.



Fig. 6.8 Layout decomposition: A single layer is decomposed in two or more masks to enhance the resolution.

Double/Multi-patterning

Although technologies keep scaling to the order of a few nanometers, lithography still uses 193*nm* wavelength light, which makes printability and manufacturability more challenging due to increased distortion. Beyond 20*nm* the use of multi-patterning is required for device fabrication. Using multi-patterning technology, a single layout is decomposed into two or more masks and manufactured through two or more exposure steps. These masks are then combined to get the original intended layout. By decomposing the layout into two or more masks as shown in Figure 6.8, the pitch size is effectively doubled thereby enhancing the

resolution [58]. To achieve this, on the design side, color (mask) assignments are used. Several techniques of multi-patterning include Litho-Etch-Litho-Etch Double Patterning (LELE DP), Spacer-is-Metal Self-Aligned Double Patterning (SIM SADP), Litho-Etch-Litho-Etch-Litho-Etch Triple Patterning (LELELE TP) and Spacer-is-Dielectric Double Patterning (SID SADP). To use these techniques the designer can include the colored masks per layer that must be multi-patterned or use a colourless flow where the foundry performs the decomposition [211].

Fin Formation

Although multipatterning brings new fabrication challenges, some of the known fabrication steps from the planar technology can be repurposed to achieve new required shapes like the 3D fins. Sidewall spacer deposition steps from planar processes are utilized to perform self-aligned double patterning (SADP). Similarly, the steps used to form Shallow trench isolation (STI) can be extended to fabricate fins by additional etching of STI areas and thereby exposing Si fins. Fins are fabricated in a regular fashion over a large area. Thereafter unwanted fins are excised and the remaining fins become a part of active areas of the devices. Hence FinFET fabrication becomes compatible with old planar CMOS processes using repurposing of existing steps, plus a few extra steps.

Shape of the Fins

Several studies have shown that FinFET performance is affected by the cross-sectional area of the fin, therefore the fin shape. Intel's 22nm node microprocessor was built with FinFET sidewalls sloping at about 8 degrees from vertical which makes more sturdy devices among other advantages [137]. Figure 6.9 shows the main types of fins analyzed in the literature. Experimental data shows that a FinFET with a rectangular cross-sectional area has better short channel effect metrics, in particular sub-threshold slope, GIDL and DIBL if



compared with a triangular or trapezoidal cross-sectional area [129]. On the other hand a triangular fin can reduce leakage current by 70% if compared with a rectangular fin [65].

Fig. 6.9 Left side: a fin with a vertical slope which presents better short channel metrics [129]. Middle: a standard fin with some degree of inclination as the one used in the 22nm Intel's node [137]. Right side: a fin with a triangular cross-sectional area that can help to reduce the leakage [65].

Middle-end-of-line (MEOL)

Middle-end-of-line (MEOL) is a new term introduced in the FinFET era. It refers to the intermediate process steps that complete the transistor formation (Front-end-of-line: FEOL) before contacts and interconnect formation (Back-end-of-line: BEOL) [22]. MEOL is necessary to provide better cell level connections with restricted patterning capabilities and multipatterning [170]. The introduction of MEOL increases the complexity of fabrication and modeling as well. For circuit designers, the added new parasitic effects from MEOL need to be considered during the design process since these parasitics have been demonstrated to be one of the dominant sources [116]. MEOL parasitics have been usually accounted at the logic gate-level parasitic extraction step using the standard EDA tools. For physical



Fig. 6.10 Different FinFET logic styles: 2-input NAND gate designs with SG and IG devices.

design engineers, the added MEOL means more complex design rules and longer debugging process, also, the layout tools must automate conformance to rules as much as possible.

6.2.5 FinFET Circuits

Since FinFET devices have much better electrostatic properties and other metrics than planar devices, new logic and wider design space exploration opportunities become available. In this section, we discuss these new changes that FinFETs have introduced at the circuit level.

Logic Styles

As discussed in Section 6.2.2, FinFETs come in two flavors – short-gated (SG) and independent-gated (IG). For IG FinFETs, the top part of the gate is etched out, resulting in two independent gates. Because the two independent gates can be controlled separately, IG-mode FinFETs offer more design styles [148, 142]. Although the gates are electrically isolated, their electrostatics are highly coupled. The threshold voltage of either of the gates can be easily influenced by applying an appropriate voltage to the other gate. Shown in Figure 6.10 is one example of different flavors of 2-input NAND gate implemented using

SG/IG gate or a hybrid of both (modified from [142]). In SG mode, FinFET gates are tied together, so they work the same as the planar devices; In IG mode, one device (with two gates) is driven by two independent signals, and some logic functions can be realized by one device; in IG-Low Power mode, one gate is disabled and acts as the reverse-biased back-gate. The designers can even mix the two types of devices and balance the tradeoff if it is allowed by the foundry. But IG gate requires one more step of etching in the fabrication step.



Fig. 6.11 ON current vs. Body bias for a 2-finger 1xnm and 7nm NMOS transistor.

Body Effect

Adaptive Body Biasing (ABB) has been used by circuit designers as an effective design technique to reduce the impact of die-to-die and within-die variations by changing the NMOS and PMOS threshold voltages independently in order to maximize performance [213]. FinFETs fabricated in bulk or SOI processes receive little benefit from controlled body effect because the channel in the FinFET is mostly in the top of the fin, away from the body. Thus



Fig. 6.12 16-input AND gate implemented with different stack height (1, 4 and 16).

the body bias techniques is not applicable to FinFET circuit design anymore [41]. To validate the above argument, we apply both reverse and forward body bias to a 2-finger transistor and simulate the ON current for both 1xnm and 7nm nodes, with the results are shown in Figure 6.11. The ON current doesn't change with the body voltage, as expected, and it indicates that FinFET devices are largely insensitive to the body effect. On one hand, this reduces the available design knobs, on another hand, this can actually mitigate the stack effect and enable higher stack height.

In some logic cells, NAND gate for example, several transistors are connected in series and stacked. In planar CMOS circuit, stack height is limited by the body effect; due to the body effect, the voltage between source and body of the top stacked transistor will increase the threshold voltage and will lead to performance degradation; if the stack height keeps increasing, the pull down current will become smaller and the circuit will become slower or might not even function correctly. For FinFET logic due to the insensitivity to the body effect as discussed above, the stack effect will be minimal and this can lead to higher stack logic cells with potential of increasing the fan-in and reducing the logic depth, thus further



Fig. 6.13 16-input AND delay simulations with different stack height (interconnect capacitance is considered).

reducing delay and leakage paths. Our first attempt of simulating a 16-input AND gate confirms the above assumption. Shown in Figure 6.12 is a 16-input AND gate implemented with different stack heights and logic depths. Figure 6.13 shows the simulated delay in 1xnm FinFET technology corresponding to different stack height. The results suggest that a stack height of 16 and a corresponding logic depth of 2 stages achieves the best performance. Another benefit of increasing the stack height is the reduction of leakage. If we assume that the leakage with stack height of 16 design is 16*I*, where *I* is the leakage of the unit-sized transistor, then the leakage for a stack height of 2 is (16+8+4+2)I, which is much larger. In summary, due to the fact that the stack effect is weak in FinFET logic, designers can increase the stack height with a relative relaxed margin to balance the tradeoffs of area, delay and leakage.

	7nm FinFET		1xnm FinFET		28nm FDSOI			130nm Bulk			textbook				
	inv	nand	nor	inv	nand	nor	inv	nand	nor	inv	nand	nor	inv	nand	nor
8	1.00	1.35	1.59	1.00	1.06	1.34	1.00	1.11	1.52	1.00	1.14	1.54	1	1.33	1.67
p	1.68	2.59	3.38	0.62	1.30	0.95	2.90	4.21	4.52	0.49	0.96	0.80	1	2	2



Table 6.2 Normalized logical effort *g* and parasitic delay *p* values

Fig. 6.14 Simulated FO4 delays for Inverter, 2-input NAND and 2-input NOR gates in different technology nodes (all values are normalized to the 7nm FO4 INV delay)

Logical Effort

The logical effort method is an approximate, simplified model to analyze the delay of a gate. The normalized delay is expressed as:

$$d = f + p = g \cdot h + p \tag{6.3}$$

where p is the parasitic delay, i.e. the delay of the gate driving no external load, and f is the effort delay, expressed as the product of logical effort g and fanout h. The logical effort g is

proportional to the complexity of a gate as a more complex gate leads to higher gate delay. The fanout h is the ratio of the output load capacitance to the input capacitance of a gate.

We extimated the g and p for an inverter, a 2-input NAND and a 2-input NOR for different technologies using simulation. For this, we use a simple simulation setup consisting of *fanout* of 1 and *fanout of 4* gate delay chains. The results obtained are summarized in Table 6.2. The values of g and p have been normalized to the respective inverter values for each technology. The table shows that the g and p values vary slightly across technologies depending on transistor sizing for different technologies. Measured normalized delays for different gates are presented in Figure 6.14 which shows that gates maintain a similar trend for increase in complexity across different technologies. NOR gates with stacked PMOS are slower than NANDs (stacked NMOS) even in FinFETs where the ratio of ON current in NMOS to PMOS is close to 1 as listed in Table 6.1.

Standard Cell Libraries

There are many tradeoffs that need to be considered when developing standard cell libraries. For example, logic offerings such as the maximum number of logical inputs on complex gates, flip-flop and latch offerings, clk buffers, drive strength for each cell and so on. As discussed in the previous sections, FinFET devices have several unique intrinsic device characteristics, and these bring several changes to the standard cell library designers. First, with planar transistors, designers can arbitrarily change transistor width in order to manage drive current. With FinFETs, due to the width quantization fact as discussed in Section 6.2.2, they can only add or subtract fins to size it and change the current. Second, since body biasing is generally ineffective, as discussed in the last section, this might lead to more logical inputs on complex gates in FinFET libraries. Coming to the physical design, the FinFET devices have periodic structures, and the optimal Wp/Wn ratio is almost 1:1, thus the FinFETs layout looks more regular, and the PMOS and NMOS regions are symmetric.

The standard cell template height (in the number of M1 wiring tracks) usually comes in several flavors. For example, a high density library might be 9 tracks tall, a high performance library might be 13 tracks tall, and a power optimized library might be 10.5 tracks tall. But in FinFET, the additional constraint of fitting a fixed number of fins within a cell complicates this [7]. Especially in most FinFET technologies, fin and metal pitches are different and have not tended to line up. Power rail connections at the top and bottom of the cell typically force the removal of 1 fin each, and typically 2 additional fin tracks must be removed in the center of the cell to accommodate gate input connections, all of these make compact FinFET cell design very complex. In addition, [7] also pointed out that to meet the multiple patterning requirement, the coloring process need to be conducted during the design of the standard cells, coloring also needs to meet density solutions (each color mask must have reasonably consistent density across the chip).

FinFET SRAM Design

SRAMs are one of the most area and power hungry components on a chip. The neverending demand for packing more functionality per area and the requirement of higher performance from processing units leads to continuous scaling of devices [30]. This scaling trickles down to smaller bitcells and enables an increase in memory array density in terms of number of bits stored per area. Hence from the density point of view, minimum sized transistors are desired in bitcells. This translates to a 1 : 1 : 1 (PU:PG:PD) fin bitcell for FinFETs (where PU is the size of the Pull-up PMOS, PD is the size of the Pull-down NMOS, and PG is the size of the pass-gate NMOS in a 6T SRAM cell). The 1:1:1 bitcell provides highest array density but it suffers from flaws in terms of lower read stability and writability [138, 30]. The constant need for voltage scaling to lower power further exacerbates SRAM readability and writability issues. This calls for alternate bitcells like the Low Voltage (LV) 1 : 1 : 2 cell and High Performance (HP) 1 : 2 : 2 cell [30] along with read and write assist techniques to improve SRAM metrics. Several assist techniques [242, 197] have been proposed and studied to improve SRAM performance and lower operational Vmin. These techniques focus on improving PD:PG strength ratio for read assists and PG:PU strength ratio for write assists. These techniques become increasingly necessary in the era of FinFET SRAM design because transistor width quantization in terms of number of fins decreases device level sizing options to improve SRAM bitcell functionality.

Thermal Effect Inversion (TEI)

Thermal behavior is one of the important device characteristics that affect the design decisions like margins, floorplan and cooling costs. It has been shown in the literature that temperature characteristics of FinFET-based circuits are fundamentally different from those of conventional bulk CMOS circuits [121, 230]. In a bulk technology, if the transistor operates in the super-threshold region, the delay increases with the temperature, and in the near/sub-threshold region, the delay decreases with the increasing temperature. While in FinFET, it has been reported that the circuits run faster at higher temperatures in all supply voltage regimes (including the super-threshold one), and this is called the Temperature Effect Inversion (TEI) phenomenon [121]. In both planar devices and FinFET devices, the threshold voltage decreases at the higher temperature, and the mobility of charge carriers in the channel decreases due to the ionized impurity and phonon scattering [198]. TEI happens due to the fact that FinFET channels are usually undoped or lightly doped, so they exhibit only a small change in mobility with temperature. It has been shown in [243] that TEI's inflection voltage approaches nominal supply and the impact of this effect can no longer be safely discounted when scaling into future FinFET and FDSOI devices with smaller feature sizes. To validate this, we simulate the delay vs. temperature for a 9-stage ring oscillator in multiple technology nodes. The simulation results are shown in Figure 6.15; the results show that for all technologies, the increased temperature slow down the devices if they work 188



Fig. 6.15 Simulated thermal characteristics (Delay vs. Temperature) in multiple technology nodes for a 9-stage ring oscillator. Blue - Super-threshold; Orange - Near-threshold; Red - Sub-threshold.

under near and sub-threshold region. Interestingly, for the 28nm FDSOI node, TEI appears across all voltages, and for 1xnm bulk FinFET node, the TEI effect has already approached 0.7V, which is only 0.1V below the nominal voltage (0.8V). Similarly, for 7nm bulk FinFET, the inversion starts from around 0.6V (0.1V below the nominal voltage of 0.7V). We can conclude that the TEI effect is indeed becoming increasingly important in current and future technologies as it will cover all of the operating voltage ranges.

The TEI effect introduces new tradeoffs and also challenges in circuit design. On one hand, a higher temperature increases the leakage and cooling budget, but, on the another hand, it helps with the performance. The benefits of TEI can be maximized with the assist of novel power management techniques that can dynamically tune the voltage or frequency based on the real-time temperature [243, 150] or novel algorithms that can determine the maximum performance under power constraints [31]. Since thermal issues also emerge as important reliability concerns throughout the system lifetime, the TEI effect can compensate some of the performance degradation introduced by reliability threats such as BTI and EM [32, 198]. The optimal operating temperature can be exploited to reduce design cost and runtime operating power for overall cooling with the proper utilization of the TEI effect.

Variability and Reliability

Reduced feature size causes statistical fluctuations in nanoscale device parameters which are known as process variations. They lead to mismatched device behaviors and degrade the yield of the entire die. In planar devices, a number of dopants must be inserted in the channel which lead to Random Doping Fluctuations (RDF) causing significant variations in threshold voltage. In FinFETs, since the channel is undoped or lightly doped, this reduces the statistical impact of RDF on V_{th} . The variability associated with line-edge roughness (LER), the random deviation of gate line edges from the intended ideal shape, which results in nonuniform channel lengths, is also lower in FinFETs. But other process variations do appear in FinFETs. Since they have small dimensions and lithographic limitations, these devices suffer physical fluctuations on gate length, fin thickness or oxide thickness [21, 17, 225]. Overall, FinFETs emerge superior to planar devices by overcoming RDF and LER, which are two major sources of process variation.

Besides process variations, which represent the time-zero process variability, timedependent variations (wearout) such as Bias Temperature Instability (BTI), Hot Carrier Injection (HCI) and Electromigration (EM) detailed in this thesis also appear to be critical for reliability considerations in FinFET era. As the technology scaling is reaching the nanoscale FinFET regime, the transistors become more susceptible to voltage stress due to the increased effective field associated with the scaling of the thin oxide. Similarly, the shrinking geometries of metal layers render higher current densities, and the tremendous number of transistors within a compact area results in higher power densities. Together, these lead to increased on-chip temperatures which potentially accelerate the wearout effects [207]. Besides, the thermal resistance (R_{th}) of the multi-gate topology and the reduced gate pitch in FinFET devices exacerbate self-heating which will accelerate aging [92]. For interconnect reliability, EM no longer can be signed off using aggressive margins, a comprehensive thermal-aware EM signoff methodology needs to be adopted for FinFET designs. New types of EM rules that are dependent on the direction of current flow, metal topology, via types, co-vertical metal overlaps etc. are required to address the potential reliability issues [1]. For FinFET wearout such as BTI and HCI, a detailed analysis will be presented in Section 6.3.

Interconnect

As the devices become smaller and smaller, the interconnect becomes more and more dominant in determining circuit performance. This is because of the yield and EM requirements, the interconnect can't scale at the same rate as the transistors. As interconnect is becoming more compact at each node below 20nm [115, 130], the interconnect RC parasitic
delay will affect the performance in a more significant way and become one of the bottlenecks on the scaling roadmap. To address this, interconnect materials such as Aluminum, Cobalt (Co) or Ruthenium (Ru) could be better alternatives due to the better sheet resistance, but there are also cost and reliability considerations in the interconnect scheme design [214]. The pitch size of the metal lines also doesn't scale down that much as the technology moves into the sub-20nm regime due to the RC parasitic and coupling consideration as well. For designers, since they don't have control over the materials and design rules, the only knob they have is the dimension of the wire. This requires to consider interconnect capacitance in the early design phase even before the physical design. The FinFET PDKs usually provide relatively accurate wire models to account this.

6.2.6 **FinFET Technology for Energy-constrained IoT Applications**

FinFETs provide improvements in power and energy consumption since they overcome the leakage problems of planar devices and deliver better performance. To further investigate this aspect, we simulate a NAND-based ring oscillator [220] across multiple technologies. The duty cycle of the ring oscillator can be tuned and in our case, it is set as 10%. Shown in Figure 6.16a is the simulated delay vs. V_{dd} , in which the values of each node are normalized to the delay at their own nominal voltages. It shows that FinFETs provide a significant performance advantage at any operating voltages, and the reduced performance due to lowering the voltage is much lower in FinFETs compared to other technology nodes as well. Figure 6.16b presents the energy vs. V_{dd} plot, similar normalization is applied. As it shows, although the minimum energy optimal points are similar for all the technologies (around 0.2 - 0.3V range), the energy of FinFET scales the best with voltage; in other words, as the voltage is scaled down, FinFETs offer more energy savings than planar devices. In Figure 6.16c, the energy delay product vs. V_{dd} is plotted. FinFETs offers the best energy efficiency for circuit operating under a wide range of voltages since, as the voltage scales down, the



Fig. 6.16 (a) Delay vs. V_{dd} ; (b) Energy/cycle vs. V_{dd} ; (c) Energy Delay Product (EDP) vs. V_{dd} and (d) Minimum EDP values across multiple technology nodes (simulated with the same NAND-based ring oscillator structure)

energy delay product doesn't change significantly for FinFETs compared to planar devices. Figure 6.16d presents the minimum energy delay product (EDP) across the four technology nodes. As technology scales, the EDP improves as expected.

The above study shows that FinFETs provide more options for performance vs. other metrics tradeoffs. For example, since FinFETs offer very good energy efficiency over a wide range of voltages, voltage scaling techniques can be very effective as IoT circuit designers

strive to maximize performance per *mW* without hurting energy. FinFET-based design will be able to support wider use of dynamic voltage frequency scaling (DVFS) and enable a wider range of applications from high-end performance critical systems to energy-constraint devices such as in IoT applications.

6.2.7 Summary - Digital Circuit Design with FinFETs

We have shown so far that FinFET devices offer significant performance improvements and power reduction compared to planar devices. Digital circuit design with FinFET broadens the design window once again. Operating voltage continues to scale down, short channel effects are reduced significantly, the process variation have been improved, the FinFET devices have lower leakage power in standby mode, etc.

Although FinFET devices offer advantages in many dimensions, they also bring challenges in the design process. FinFET devices have non-standard shapes and require complex modeling of the parasitics in the TCAD tools. Moreover, the physical layout-dependent effects have a significant impact on the metrics. Therefore, the design tools and design flows need to be able to assist the designers to build circuits that accurately correlate to the models. During the design process, extraction plays a big role to obtain accurate timing analysis and power estimation for FinFETs, so enhancement to the foundational EDA tools, in particular SPICE simulations, extraction and physical verification that operate on part of the design below the first metal layer are required [101]. Interconnect resistance is becoming more important, so IR drop and power-grid design becomes more critical. Besides, to meet the double/multipatterning requirements, the standard cell, floorplanning, placement and route (P&R) need to be colored correctly. For example, during power planning, all power rails need to be free of double patterning violations. Similarly, all the placement of standard cells and hard macros need to be double patterning-compliant. Physical verification (e.g. DRC) engines need to be able to check and guide the designers to meet the double patterning rules. More verifications are required, and more checkpoints need to be inserted during the design phase to make sure the design specification is met.

For custom designers and standard cell designers, all of the blocks require a redesign due to the following reasons. First, the options of sizing are less granular due to the width quantization fact in FinFET, getting more drive strength will require more fins in parallel. Second, the thermal behavior and options available to circuit designers are different than what they may be used to with planar devices. For example, body biasing will be impractical, thermal effect inversion (TEI) fact introduces new tradeoffs, higher fan-in and complex logic are possible due to the insensitivity to the stack effect. As dozens of new and complicated design rules arise for FinFET devices, physical design efforts are increasing, but the bright side for FinFET devices is the more regular layout and equal P and N regions, and because of this, the foundry usually provides a template layout on which fingers and gates are already placed, physical designers don't need to start from scratch, but the layout tools still need to automate conformance to rules as much as possible.

FinFETs also offer more design options for trading performance with other metrics. As discussed in Section 6.2.6, one major design optimization benefit of FinFETs is much higher performance with the same energy budget. Similarly, they consume much lower power and energy to achieve equal performance to planar devices. This essentially gives designers the ability to extract the highest performance for the lowest power, which is a critical optimization for battery-powered devices. Since FinFETs have lower leakage and can operate faster, the circuit can afford to have more and fine-grained power gating structures to further save power in standby mode. Runtime techniques like DVFS can be used with a lower cost to maximize energy efficiency. On top of all these benefits, the circuit can operate in near-threshold to save energy with lower performance penalties [161].

As more transistors fit on one chip in the FinFET era, the design flow needs to be able to handle big designs which have billions of transistor at a fullchip level, thus optimizing the runtime and reducing peak memory are necessary, and more parallelism is required. Because of the increased complexity and number of instances on chip, an increasing number of signoff corners are required to cover process and environmental variations. Addressing these new challenges together with the new, more complex design-for-manufacturing rules, including double/multi-patterning, along with the increasing design scale, require close collaboration between the foundry, tool vendors and designers to fully take advantages of what FinFETs have to offer.

6.3 When "things" get older – Exploring Transistor Aging in IoT Applications

6.3.1 Motivation

In previous sections, we looked into opportunities and challenges FinFET technology has introduced. Due to the great energy efficiency and highest level of integration, FinFET is going to play an important role in many emerging applications such as Internet of Things (IoT) [227, 187]. As the number of connected devices and connections between human and "things" increases rapidly, Internet of Everything (IoE) emerges as a wider concept of connectivity platform from the perspective of modern connectivity technology use cases, where "things" are key components and lay the foundations for the massive interactions with the world [91, 95]. IoT is a general-purpose technology by nature ranging from health care, transportation to agriculture and almost all aspects of life [9]. These applications impose common requirements for IoT devices like they should have small form factors in terms of physical dimensions and weights. Since most of these devices are battery-powered or batteryless, they require high energy efficiency and extreme low power consumptions. In addition to these characteristics, IoT devices need to withstand hostile environments such as increased and highly variable temperatures and voltage noise [9]. More importantly, IoT

nodes are required to operate reliably for a long lifetime (e.g. decades), which translates to reliability challenges, especially device degradation-induced circuit aging⁴. As has been discussed in details in previous chapters, on-chip elements such as transistors and metal wires age gradually when under use, and this can lead to potential permanent failures. Many of the IoT applications (such as automotive or implantable medical devices) require almost zero error during the whole lifetime. Harsh environment such as high temperature accelerate aging. As number of on-chip elements scales up, more transistors are susceptible to aging and this leads to the increase of the system failure rate. These advanced technology nodes impose more aging issues than previous generations due to self-heating, reduced oxide thickness, narrower metal and increased current density [202].

Since IoT is a wide concept and circuit aging is a threat to IoT lifetime, it is necessary to understand how current and future IoT systems are impacted by aging and how to deal with it in this context. In this section, we will look into these aspects by conducting extensive circuit-level simulations with foundry-calibrated aging models in advanced FinFET node. As aging is highly dependent on application behaviors that define the operating voltage, temperature and active time, we perform a survey of existing IoT applications and classify them based on aging-related metrics. The aging behaviors in each category will be discussed in the following sections.

6.3.2 Previous Work on Aging in IoT Domain

As has been detailed in Chapter 2 and 3, the primary cause of aging is the electrical stress across the on chip components such as transistor, dielectrics and interconnects. In general, there are mainly three dominant causes of aging in semiconductor devices. Bias temperature instability (BTI) and Hot carrier injection (HCI) lead to transistor degradations [183] and Electromigration (EM) leads to metal wire resistance increase. The main mechanism of BTI

⁴This thesis focuses on circuit aging. Battery aging and socket (and holder solder) aging are out of the scope of discussion.



Fig. 6.17 Transistor Aging: HCI occurs mainly during switching; PBTI happens when NMOS is under stress; NBTI happens when PMOS is under stress. BTI aging partially recovers during OFF states.

and EM have been discussed in previous chapters. HCI shares many similarities to BTI, both of them impact transistor parameters (e.g. threshold voltage V_{th} and carrier mobility μ) at a level that depends on the operating environment and usage of the circuit. As illustrated in Figure 6.17, BTI is mainly caused by constant electric fields degrading the dielectric. HCI is also caused by electrical field, but it happens mainly on drain side and primarily occurs during switching. While BTI is partially reversible HCI is an irreversible effect. The parameter shift is also highly temperature dependent since temperature affects the interface trap generation. In this chapter, we mainly focus on transistor aging (BTI and HCI) because the current density in IoT chips is relatively low, thus EM is less impacted compared to other aging mechanisms.

Transistor aging has been explored for a long time in aerospace and safety applications but it didn't gain interests in consumer devices until very recently. For example, in automotive or industrial IoT applications aging happens even when the system is inactive most of the lifetime due to the continuous constant stress across transistors. These devices need to function under all possible scenarios during their expected lifetime [147]. For example, some medical implants will require a reliable operation for more than 50 years [61]. While a car today may sit idle 90% to 95% of the time, an autonomous vehicle might only be idle 5% to 10% of the time [54]. Previously a lot of attention has been paid on battery or package aging for IoT applications [181] while a few studies looked into circuit aging in this domain. [122] introduced a method to obtain multi-threaded switching activity signatures for aging analysis in IoT applications but the focus was on the architectural level framework. Similarly, [196] proposed a solution that leverages the workload dependent reliability analysis for early product failure rate calculations for automotive applications. [61] provided a device level experimental study of BTI aging in ultra-low power applications. [215] proposed a unified model which captures the joint impact of RTN, BTI and PV within a probabilistic reliability estimation for NTV circuits. Most of the previous studies focused on circuit aging in a very specific application or framework. Also previous works used analytical aging model which can lead to inaccurate predictions. The contributions of this part of work are:

- We study aging impact with foundry-calibrated model instead of predictive models.
- We investigate transistor aging on a wide spectrum of IoT applications and provide a deep and realistic understanding of how aging affects each IoT category.
- Several solutions for addressing aging in IoT including active recovery are also discussed.

6.3.3 IoT Application Domains

The hardware requirements of the IoT devices are determined by how and where they are deployed [9]. To develop quantitative understanding of these requirements for IoT nodes, we surveyed published SoCs and commercially available IoT products, ranging from agriculture/environmental sensors, automotive, industrial processes to medical implantables,

$\bigtriangleup:$ Specification for current commercial	product	s (in ma	arket); ⊀	: Speci	fication	for futu	ire prod	uct (und	der rese	arch)		
Applications	Tei	nperati	ure	Co	re Volt:	age	Rec	Lifetime	e ents	Ac	ctive Tin	ne
	Г	Μ	Η	Γ	Μ	H	Г	Z	H	Г	Μ	Η
1 - Implantable/Heathcare		\triangleleft		*	*	\triangleleft			\triangleleft			\triangleleft
2 - Consumer Electronics/Wearables	\triangleleft	\triangleleft		*	*	\triangleleft	\triangleleft				\triangleleft	
3 - Automotive			\triangleleft			\triangleleft			\triangleleft		\triangleleft	
4 - Industrial Processes	\triangleleft	\triangleleft	\triangleleft			\triangleleft		\triangleleft				\triangleleft
5 - Public Transportation	\triangleleft	\triangleleft				\triangleleft		\triangleleft			\triangleleft	
6 - Energy Management	\triangleleft	\triangleleft	\triangleleft			\triangleleft		\triangleleft			\triangleleft	
7 - Smart Homes/Buildings/Cities	\triangleleft	\triangleleft		*	*	\triangleleft		\triangleleft				\triangleleft
8 - Retailing/Malls		\triangleleft				\triangleleft		\triangleleft			\triangleleft	
9 - Agriculture/Environmental	\triangleleft	\triangleleft				\triangleleft	\triangleleft				\triangleleft	
10 - Wildlife/Nature Preservation	\triangleleft	\bigtriangledown				\bigtriangledown			\bigtriangledown		\bigtriangledown	
- For temperature, L - Low (${\leq}27^{\circ}C),M$ -	- Mediu	m (27°C	$C - 100^{\circ}$	°C), H -	High (>100°C	;(;					
- For voltage, L - Low (Sub-threshold), N	A - Med	ium (N	ear-thre	shold),	H - Hig	h (Nom	inal Vo	ltage);				
- HOr litetime requirement 1 OUV (< 4 VP9)				10	VEARCI	Ĭ		101691				

Table 6.3 Summary of IoT Applications Specifications (Aging-related Metrics)

- Here, "active" means transistors are under stress. In many IoT applications, even when the systems are in "sleep" mode, many For intentite requirement, $Low (\geq 3$ years), M - intention (3 years - 10 years), H - High (> 10 years);

circuit blocks (such as accelerometers in a wristband) are still ON and under aging stress. L - Low (Active < 20% of the lifetime), M -Medium (Active 20% - 80% of the lifetime), H - High (Active > 80% of the lifetime). smart cities and consumer electronics. The results are summarized in Table 6.3 and are discussed below.

We classify the existing IoT applications into ten groups mainly based on the usage and scale of users [9, 72]. We mainly study aging-related metrics, i.e. voltage, temperature, lifetime requirements and active time which is how long transistors are stressed. IoT chip in all categories are found to operate at super-threshold voltages during active phases of computation for speed purposes even in battery-operated systems. Most of commercial low-power IoT architectures achieve energy efficiency through heavily optimized deep sleep modes or minimization of the unnecessary on-chip components [9]. There is a lot of ongoing research on operating IoT systems completely in near/sub-threshold region to achieve major energy efficiency improvements, especially in applications such as medical devices, sensors and wearables where energy harvesters can be adapted (application 1, 2, 7 in the table). But this comes at the expense of performance and increased sensitivity with respect to variations. In this work, we focus on IoT chips that operate at nominal voltage only. Application 1 and 2 represent personal IoT where implantable devices usually operate continuously at human body temperature while consumer electronics such as wearables are exposed to environment. Implantables are always active and require a long lifetime (almost human being lifetime) while wearables have a relatively shorter life cycles (around 3 years) and are inactive most of the lifetime. Applications 3 and 4 are industrial IoT in which automotive sensors monitor the state of the vehicle and mostly reside inside engines. They operate at very high temperature and require a reliable operation throughout car's lifespan of more than 10 years. Similar monitoring strategies are used in industrial environment such as storage warehouse or product lines. IoT devices also enable ubiquitous sensing in city and home infrastructures (applications 5 - 8) and are installed both indoors and outdoors. Thus they experience room or environmental temperatures and have relatively strict lifetime requirement since frequent checking and repairs are not practical. Applications 9 and 10

represent environmental IoT applications. They have similar temperature requirement as city/home-scale IoT. The agriculture sensors usually last for one cycle of crops but other environmental IoTs need to last longer because they are distributed at a very large scale and many of them are not quite accessible physically once they are installed.

6.3.4 Simulation Results

In general, process, voltage and temperature (PVT) variations impose additional timing margins that stretch the clock cycle. Aging impacts circuit in a similar way. Aging-induced performance loss requires guardband for margining. In this section, we study impact of aging with the foundry-calibrated models in FinFET technology using Cadence reliability simulator RelXpert (integrated in Virtuoso ADE). Both BTI (including recovery) and HCI mechanisms are captured in this model.



Fig. 6.18 Single transistor ON current degradation due to aging under **DC nominal** voltage stress.



Fig. 6.19 Single transistor ON current degradation due to aging under AC nominal voltage stress with 50% duty cycle. Degradation is about half of the DC stress case due to recovery.

Single Transistor Aging under Different Conditions

Transistor aging causes parameter shift, such as increased threshold voltage V_{th} and reduced mobility μ and this leads to reduced current. Two sets of simulations (DC and AC stress) are run at different temperature and lifetime conditions for single transistor under nominal voltage. The ON current degradation (%) is plotted. Figure 6.18 shows the DC stress case, in which the PMOS transistor ages continuously during the lifetime without any recovery giving the worst case estimation of NBTI aging. While Figure 6.19 shows the results where the transistor is under AC stress with 50% duty cycle allowing BTI recovery period following stress. This includes both BTI and HCI degradations and provides an average aging estimation. As a reference, Monte Carlo simulations show that σ/μ for ON current is around 7%, which indicates the design margin with respect to process variation. The aging-induced degradations are comparable to this, and in high temperature cases, it can be more than 10%. An observation from Figure 6.18 and 6.19 is that AC stress-induced

degradation is almost half of the DC case and indicates that the degradation under the same temperature and lifetime condition is almost linearly proportional to the stress time. This assumption will be used in the following sections for estimating how the active time of each application affects aging.



Fig. 6.20 Simulation setup (an example of datapath): Aging can lead to timing failures such as setup violation by slowing down the datapath. Designers should take extra margins based on aging impact.

Aging-induced Timing Failures in IoT Scenarios

Single transistor aging causes the degradation of ON current which will further impose timing error at the circuit level and failures at the system level. Shown in Figure 6.20 is a typical data path from the output (Q) of the launch flop to the data input (D) of the capture flop. Aging slows down each unit ($t_{datapath}$ and t_{setup} increase) and this could cause setup time violations. This effect becomes more significant in data paths which have large logic depth. Extra timing margins need to be added to meet the setup timing requirement. To quantitatively study this margin under different operating conditions, we simulate a similar data path consisting of a combination of inverters and buffers. The margin is found by increasing the clock period until the launched data is correctly captured. Figure 6.21 plotted the necessary margin vs. temperature and active time (how long the transistor is stressed)



Fig. 6.21 Normalized Timing Margin vs. Temperature and Active Time: Margin shown on Y-axis is normalized to the required aging margin for datapath (shown in Figure 6.20) for 2 years at room temp (27°C).

at nominal voltage. The result has been normalized to the necessary aging-induced timing margin at 27°C with an active time of 2 years. This baseline margin is also equal to margin for temperature variations from 27°C to 110°C at time zero.

Based on the simulated results in Figure 6.21, we map the operating conditions of different IoT application listed in Table 6.3 and list the estimated aging margin for each category in Figure 6.22. As shown in the table, even in the same application category, the IoT devices



Fig. 6.22 Estimated aging margin for different IoT applications: X-axis corresponds to IoT application index in Table 6.3, Y-axis shows the normalized design margin and the error bars show design margin range within each category.

operate at different temperatures and active time. The error bar in the figure provides an estimated range of the necessary margin which also indicates the estimated aging levels for each application. The top two aging-critical applications are 4 and 6 which correspond to industrial processes IoT and energy management IoT where high temperature and long active time are expected. These applications impose a more than $10 \times$ design margin compared to



Fig. 6.23 6T SRAM read current degradation with aging for different temperatures (nominal voltage).

the baseline margin. In the second tier, automotive IoTs (application 3) suffer huge aging due to high operating temperature. Most of the city scale and environmental IoTs (applications 5 - 10) operate under environment temperature but have a relatively long active time and hence they lie in the third tier ($3 \times to 5 \times$). Similarly, implantable devices operating at body temperature need to operate reliably for a long life span, so they are also on the third most critical aging level. As consumer electronics or wearables (Application 2) are usually updated within timescale of 2 years, so they are the least aging-critical, but even so, they need to be margined properly for aging to guarantee the reliable operations spanning their lifetime.

Aging Impact on IoT SRAMs

SRAMs act as an external cache for many of today's IoT applications. They usually occupy the largest chunk of SoC and may interact with multiple cores. Hence it becomes imperative to study the impact of aging on SRAMs because the access-time and drive strength degradation may lead to timing failures across the chip. The access-time is directly proportional to the SRAM read current, I_{read} . Figure 6.23 shows the I_{read} degradation of a 6T SRAM across different temperatures for different active time. The I_{read} values have been normalized to time-zero I_{read} at 27°C. It shows that many applications will incur more than 5% degradation during their lifetime while some critical applications such as Industrial IoT facing more than 20% degradation in read current. Such a huge loss in SRAM performance due to aging will potentially lead to fatal timing errors and hence should be taken care of during design process. The design process should also aim to appropriately assign timing margin for aging based on target applications.

6.3.5 IoT Lifetime: Battery vs. Chip Lifetime

As battery replacement is not an option both due to the large numbers and inaccessibility of nodes in many IoT applications, the foremost requirement is that they can't rely on constant battery change. Thus the most common ways of defining lifetime of a battery-powered IoT system is by battery lifetime which is the time a node will operate in its normal mode without replacing the battery [9]. It is given by

$$T_{lifetime|battery} \sim E_{battery} / P_{average}$$
 (6.4)

As transistor aging could potentially lead to chip failure that might not be recoverable as discussed in the last section, the aging-induced chip lifetime should also be considered to



Fig. 6.24 IoT lifetime: Chip lifetime and Battery lifetime depend on different factors, but they can affect each other indirectly. Two lifetimes together determine the lifetime target of an IoT application.

determine the final IoT system lifetime, which is given by

$$T_{lifetime|Final} = min\{T_{lifetime|battery}, T_{lifetime|chip}\}$$
(6.5)

$$T_{lifetime|chip} \sim F(voltage, temperature, active time)$$
 (6.6)

Although battery lifetime and chip lifetime depend on different factors, they also impact each other indirectly. Illustrated in Figure 6.24 is a suggested flow for closing the IoT lifetime loop as part of the design cycle. The system lifetime target is defined by the applications and specifications, which also constrain the battery size, weight and type. On the right branch, the battery lifetime is determined based on Equation 6.4. Design knobs such as voltage, power modes and active time can be tuned to achieve lower power consumption while fulfilling the performance requirement. Some of these design knobs are also limited by application itself e.g. implantable devices (Application 1) need to be active continuously and require fast response. On the left branch of the IoT lifetime loop, chip lifetime affected by aging is also highly dependent on knobs such as voltage, temperature and active time. To guarantee that chip lifetime is longer than or equal to the battery lifetime and the overall lifetime target; the design margin needs to be reserved. But as shown in the previous section, this margin can be very large and can translate into wasted energy in the early lifetime, which in turn shortens the battery lifetime. Alternatively, aging can be addressed in run time by adaptive solutions to reduce the design margin, while the additional sensors and circuitry will add to the power budget. To leverage these tradeoffs and meet the expected lifetime, careful design decisions considering both chip and battery lifetime are required. More details are discussed in the following section.

6.3.6 Potential Solutions for IoT Circuit Aging

Adding design margin is currently the most common way of addressing aging in the design flow. This is a static solution where all transistors are margin-degraded to a certain amount based on the operating conditions. The difference in performance of aged cell versus the original cell is computed and the ratio (aged/fresh) is used to derate cells in the design. But the large margin in some applications (shown in Section 6.3.4) can be very conservative and introduce performance penalty in the early lifetime. An alternative solution would be to either recover aging using the on-chip solutions this thesis has proposed in Chapter 4 and 5 or adapt to it dynamically so that the design margin requirement can be potentially relaxed. This section will briefly discuss several such solutions for IoT applications.

Lowering the Operating Voltage

Scaling the system operating voltage down into near/sub-threshold regime has been known to be a very effective way of reducing the energy per computation and extending the

battery lifetime especially in health care and body sensor IoT applications [9]. Meanwhile, aging has an almost power law dependence on stress voltage [183]. Operation at lower voltages suppresses aging significantly. But the challenges are performance loss and increased sensitivity with respect to variations. The easy solution for performance degradations is to raise the operating voltage back again when necessary to meet the speed requirement, but this in turn will accelerate aging on the entire chip. One potential approach could be to have fine-grain voltage domains which can ensure that voltage boosting is kept small enough so that aging does not introduce large degradations and the impact is only constrained to certain sub-block. Fine-grain voltage domain can also maximize the opportunities to correct variations in paths that are critical due to variations. But this approach certainly leads to significant area overhead and design effort. These tradeoffs need to be leveraged based on the budgets that are defined by IoT applications.

IoT Circadian Rhythms: Active Recovery during Standby

Another strategy for energy savings in IoT circuit design is to put the device in "standby" state as long as possible. This is feasible since the devices don't need to be active all the time in many applications shown in Table 6.3. As we have shown in this thesis that aging mechanism such as BTI is recoverable when the transistor is OFF [183]. Hence these standby periods can be utilized for recovery. Figure 6.25 illustrates the power and aging profile of a typical IoT node, where the sensing activity is periodic and triggered by some form of real-time events. One solution to save power and reduce aging is to operate the whole processing unit at the lowest voltage level while maintaining its state in retention mode. Although there is some recovery when switching from high voltage to lower voltage but aging will still accumulate since transistors are still under stress (at a lower level). An alternative solution would be to turn off certain blocks completely such as fixed function units by power gating and save the state in retention registers. This will result in partial recovery because



Fig. 6.25 Power and Aging profile of a typical IoT node: This figure is for illustration only, the height and width are conceptually marked. For aging profile, Y-axis "Aging" corresponds to aging-induced metric change such as ΔV_{th} or timing margin (reduced performance).

BTI has a component which cannot be recovered just by switching off the transistors [68]. The third option would be to use the wearout-aware power gating structure described in Chapter 4.2.3 to reverse gate bias the transistor and heal it faster. This approach also helps to recover BTI aging components which are not recoverable at zero bias. This approach enables maximum recovery during standby mode. Although the last two approaches come with additional power and area overheads for the power switches, logic retention, signal isolation and additional floorplanning constraints [9]. The third solution also introduces one more



Fig. 6.26 Conceptual illustration of dynamic margins to enable one chip across multiple IoT applications (BOL - Beginning of lifetime, EOL - End of lifetime).

voltage source and domain (overhead is listed in Table 4.2). But for extremely aging-critical applications such as application 3, 4 and 6 shown in Figure 6.22, these overheads are justified to prevent system failures.

Dynamic Margins across Multiple IoT Applications

To minimize design costs, circuit designers and chip vendors usually try to use one SoC design across multiple IoT applications. Even within one IoT application category, the operating conditions may change. These variations necessitate run-time compensation of aging such as techniques proposed in [46], where aging events (e.g. delay change) are tracked over operating periods. Once the failure flag (e.g. timing failure) is triggered, adaptive



Fig. 6.27 Chapter 6 Highlights.

solutions such as dynamic voltage and frequency scaling (DVFS) or error correction are employed to compensate the degradations. But pure dynamic solution can be very limited and costly due to the limited tunability of metrics. A combination of static and dynamic margin methods is a more optimal approach. Figure 6.26 illustrates a potential solution where targeted IoT applications are binned based on estimated aging levels. The static margin can be added based on the lowest level in the group, dynamic solutions are also applied for adapting to the worst-case operating conditions. Compared to the purely flat guard-band based approach, the combined static and dynamic margining solution is able to leverage the power-aging tradeoffs while adapting to a wide range of IoT applications.

6.4 Conclusions

It has been almost a decade since FinFET devices were introduced to full production, FinFETs present a new frontier for the electronics industry and have enabled high performance applications such as supercomputers. While energy-constrained applications such as IoT industry is in the process of updating the technology node to FinFETs. As highlighted in Figure 6.27, in this chapter, we first studied the changes since the advent of the FinFET devices and addressed the challenges we face with these devices. FinFETs offer benefits in many dimensions such as the significantly improved power and performance metrics and lesser short-channel effect. FinFETs endeavour to offer advantages of future scaled devices while offsetting the problems introduced by many generations of planar CMOS scaling. But new challenges also appear due to many unique properties which FinFETs have shown. This part of work attempts to add to the growing FinFET design knowledge base.

The diversity of IoT applications and markets leads to a plethora of requirements for IoT reliability. The second part of this chapter showed that FinFET transistor aging introduces new challenges for IoT domain. Our study demonstrated that transistor aging should be carefully addressed earlier in a system design cycle. We also presented a set of static and dynamic solutions (e.g. active recovery) to compensate and fix aging in IoT systems.

The FinFET study work presented in this chapter has been published in [J1], and the IoT aging work has been published in [C3] and [P2].

6.5 Acknowledgements

The work presented in this chapter was supported by an IEEE CASS Pre-doctoral Fellowship, by NSF grants CCF 1619127 and CCF 1543837, by DARPA under the UPSIDE and PERFECT programs and by the Center for Future Architecture Research (C-FAR), one of six SRC STARnet Centers, sponsored by MARCO and DARPA. I'd like thank my fellow labmates Vaibhav Verma, Patricia Gonzalez-Guerrero and Sergiu Mosanu for discussions and helping with simulations.

Chapter 7

Conclusions and Future Directions

7.1 Summary of Contributions

The primary goal of this thesis has been to explore an effective recovery solution to completely repair both BTI and EM wearout effects. In this respect, this thesis provided experimental demonstrations, on-chip implementations and design methodologies for the accelerated self-healing techniques, which bring a new dimension for mitigating wearout issues in an effective way (as shown in Figure 7.1). Since accelerated self-healing is orthogonal to other wearout mitigation techniques and is feasible to implement with a relatively low cost, it can potentially be integrated with other techniques to further lower the cost and leverage the tradeoffs. This thesis aims to build the infrastructures for these directions by providing experimental evidences, circuit IP blocks and tradeoff analysis. We also looked into wearout effects in advanced FinFET nodes in emerging applications such as IoT. In summary, this thesis contributes to the reliable system design with the following accomplishments:

• Performed intensive hardware measurements on FPGAs and on-chip metal lines to understand recovery behaviors for BTI and EM wearout. The measurement results and



Fig. 7.1 Illustration of Accelerated Self-Healing as a new dimension for mitigating wearout effects.

conclusions can be used by modeling community to develop accurate recovery models for both wearout effects (Chapter 2 and 3).

- Developed the BTI gate level analytic models that can be used to predict recovery rates, and can be used together with higher level models to project system resilience (Chapter 2).
- Demonstrated that both BTI and EM recovery can be made active, and the irreversible components can be completely avoided. These properties can translate into huge benefits such as reduced design margin, less tracking power overhead and less area. Recovery can potentially lead to the wearout-aware and accelerated self-healing system (Chapter 2 and 3).
- Implemented and instrumented accelerated self-healing on chip by designing a full set
 of circuit blocks that are able to activate and accelerate both BTI and EM recovery.
 A compact circuit scheme for assisting both BTI and EM recovery and supporting
 multiple recovery modes is also designed (Chapter 4).

- Several novel types of wearout sensors are designed. Two of them are for tracking BTI wearout, in which one is designed for separating PBTI and NBTI, another is designed to monitor BTI wearout and recovery. The designs can be incorporated into standard synthesis flow (Chapter 4).
- Provided a series of potential cross-layer implementations especially at the architecture and system level. These solutions are ale to utilize some of the intrinsic sleep behaviors for self-healing. Recovery-driven design methodology is explored to enable such implementations that lead to a full accelerated self-healing system. Potential overheads are also commented (Chapter 5).
- Performed a comprehensive study across multiple technology nodes and identified the design challenges for FinFET digital circuit. Key findings included thermal effect inversion, short-channel effects, logic efforts, variations and reliability, body effect and more. This study can serve as an educational material and contribute to the growing knowledge base and design experiences for designers who are adapting to the new technologies (Chapter 6).
- Performed a first-ever complete study on impact of circuit wearout in IoT applications. The results indicated that wearout issues are very critical in several IoT domains. It is an important factor for determine the overall IoT system lifetime. This study can potentially guide IoT circuit and system designers on making high-level decisions on whether wearout effect need to be addressed for the targeting applications (Chapter 6).

7.2 Future Directions

In the next decade, CMOS technology is still predicted to be the most robust and cost effective solution for designing chips. Technology scaling, although being slow, still continues in some extent. Thus we arout issues will become more and more pronounced

together with the increasing demands for robust operations within an extended lifetime by many emerging applications such as autonomous driving, medical health and robots. Moving forward, there are many directions can benefit or be inspired from the research results of this thesis to ensure a reliable system. In this section, we will list a few of such candidates.

7.2.1 Accelerated Self-Healing in Emerging Technologies

This thesis mainly looked into recovery behaviors in CMOS circuits and interconnects. There have been increasing interests on other non-CMOS emerging technologies, specifically for non-volatile memories (NVM). Examples of such technologies can be STT-MRAM (Spintransfer torque), resistive RAM (RRAM) and phase-change RAM (PCRAM) that promise high performance, low power consumption, and unlimited endurance. These devices might exhibit their own versions of wearout, but it is likely that the fundamental similarities with CMOS wearout mechanisms still hold. For example, for STT-MRAM, wearout is due to the repeated tunneling through the magnetic tunnel junction (MTJ) with the large current during write. This is very similar to BTI where current leads to the charge trapping. The notion of "accelerated self-healing" is highly possible to be applicable to emerging technologies, this is certainly a direction that can be further explored. As it is still debating which memory hierarchy these technologies are the best for, there is good opportunity to consider wearout and recovery in the loop and make the optimal decision to ensure reliability of the systems.

7.2.2 Exploring Other Sources for Accelerating Recovery

The fundamentals of accelerated and active recovery are that the external techniques (such as high temperature and reverse bias demonstrated in this thesis) somehow affect the energy levels of the atoms or charge carriers. Ultraviolet (UV) light has long been used for erasing EPROM (Erasable Programmable Read-only Memory) which can be electrically programmed but can't be electrically erased. The charge mechanism here is the photoelectric effect where the photons directly hit and scatter the electrons from the floating gate. Therefore UV recovery has a very high potential of removing most of the trapped charges and completely rejuvenating the flash memory or CMOS circuits. Moreover, typical EPROM erasure takes only about 30 minutes to completely remove. Therefore we expect the time to remove the recoverable trapped charges to be of the same order of magnitude. UV recovery can be comparable to thermal recovery.

7.2.3 Integrating Wearout and Recovery in EDA Design Flow

Most of the current EDA tools and flows focus on optimizing timing, power and area. There are very few work on optimizing the design for robustness. Most of the time, wearout is addressed still by guardbanding, which can lead to overestimation. As we have demonstrated that recovery is very effective in this thesis, it is promising to utilize some of these behaviors and design for resilience. The very first step is to develop the device level models that capture all the recovery behaviors observed in this thesis, and these models can be used for circuit simulation, or can be instrumented in cell libraries so that they carry the wearout and recovery information through the whole flow. This thesis can also serve as the experimental evidence for validating the models or design flows.

7.2.4 Dynamic Wearout Management by Self-learning

As we expect that wearout and thermal sensors will be distributed across the whole chip, an interesting direction is to develop learning algorithms that are able to predict the circadian rhythms based on the history data collected from the sensors. Instead of instrumenting proactive recovery with fixed periods, these learning algorithms can guide the scheduler during run time. Another direction is to explore how to adjust the dynamic thermal management policies to enable thermal recovery without hurting other metrics. In the past, most of these policies focused on reducing temperature to alleviate wearout, there hasn't been efforts on recovery aspects yet. So there are good opportunities to combine dynamic thermal and wearout management techniques for enabling recovery in a smarter way.

7.2.5 Teaching Wearout and Recovery as Part of the VLSI Classes

There are very few materials on CMOS wearout and recovery in a conventional VLSI textbook. Part of the reason was because this topic is very new and is still under active research. But as reliability becomes as important as power, performance and area metrics, it is important to teach some aspects of wearout (e.g. tradeoffs, design techniques, recovery aspects) for students to understand the basic tradeoffs. Since this thesis has produced lots of new experimental results based on modern computing platforms such as FPGAs and processors, and we also looked into wearout issues in emerging applications like IoT with the advanced nodes such as FinFET, so some additional efforts can be added to make part of this thesis educational materials to transfer the knowledge.

Appendix A

List of Publications

A.1 Peer-Reviewed Journals

- [J1] X. Guo, V. Verma, P. Guerrero, S. Mosanu, M. Stan, "Back to the Future: Digital Circuit Design in the FinFET Era," *Journal of Low Power Electronics* (JOLPE), Vol. 13, No. 3, pp. 338–355, 2017. (Invited)
- [J2] X. Guo, M. Stan, "Implications of Accelerated Self-Healing as a Key Design Knob for Cross-Layer Resilience," *Integration, the VLSI Journal*, 56 (2017): 167-180, Elsevier.
- [J3] M. El-Hadedy, X. Guo, M. Margala, M. Stan, and K. Skadron. "Dual-Data Rate Transpose-Memory Architecture Improves the Performance, Power and Area of Signal-Processing Systems." *Journal of Signal Processing Systems* (JSPS), Vol. 88, No. 2, pp. 167-184, 2017.
- [J4] M. El-Hadedy, X. Guo, M. Stan, K. Skadron, "Crypt-Pi: A Light and Fast Crypto-Processor for IoT Applications," *To be submitted*.
- [J5] K. Mazumdar, X. Guo, R. Zhang, M. Stan, "Charge-Recycled Power-Regulation with Stacked Loads and Stacked Switched-Capacitors," *IEEE Design and Test of Computers* (D&T), submitted.

A.2 Peer-Reviewed Conferences and Workshops

- [C1] P. Guerrero, X. Guo, M. Stan, "SC-SD: Towards Low Power Stochastic Computing using Sigma Delta Streams," *submitted*.
- [C2] A. Roelke, X. Guo, M. Stan, "OldSpot: A Pre-RTL Model for Aging and Lifetime Optimization," *submitted*.

- [C3] X. Guo, V. Verma, P. Guerrero and M. Stan, "When "things" get older Exploring Circuit Aging in IoT Applications," In Proc. of International Symposium on Quality Electronic Design (ISQED), Santa Clara, CA, 2018.
- [C4] D. Kamakshi, X. Guo, H. Patel, M. Stan and B. Calhoun, "A Post-Silicon Hold Time Closure Technique using Data-Path Tunable-Buffers for Variation-Tolerance in Subthreshold Designs," In *Proc. International Symposium on Quality Electronic Design* (ISQED), Santa Clara, CA, 2018.
- [C5] S. Eldridge, V. Verma, X. Guo, A. Roelke, K. Swaminathan, N. Chandramoorthy, M. Cochet, A. Buyuktosunoglu, C. Vezyrtzis, R. Joshi, M. Ziegler, M. Stan, P. Bose, "VELOUR - Very Low Voltage Operation Under Resilience Constraints," In Proc. of The Government Microcircuit Applications and Critical Technology Conference (GOMACTech), Miami, FL, 2018.
- [C6] M. El-Hadedy, X. Guo, X. Huang, M. Margala, "RE-HASE: Regular-Expressions Hardware Synthesis Engine," In Proc. of The 3rd International Workshop on Heterogeneous High-performance Reconfigurable Computing (H2RC'17), in conjunction with SC, Denver, CO, 2017.
- [C7] X. Guo, M. Stan, "Deep Healing: Ease the BTI and EM Wearout Crisis by Activating Recovery," In Proc. of IEEE/IFIP International Conference on Dependable Systems and Networks Workshop (DSN), Denver, CO, 2017.
- [C8] M. El-Hadedy, X. Guo, M. Stan, K. Skadron, "PPE-ARX: Area- and Power-Efficient VLIW Programmable Processing Element for IoT Crypto-Systems," In Proc. of NASA/ESA Conference on Adaptive Hardware and Systems (AHS), Pasadena, CA, 2017.
- [C9] S. Eldridge, K. Swaminathan, N. Chandramoorty, A. Buyuktosunoglu, A. Roelke, X. Guo, V. Verma, R. Joshi, M. Stan, P. Bose, "A low voltage RISC-V heterogeneous System: boosted SRAMs, machine learning and fault injection on VELVET-bonus," In *Proc. of Workshop on Computer Architecture Research with RISC-V* (CARRV), in conjunction with IEEE MICRO, Boston, MA, 2017.
- [C10] X. Guo, M. Stan, "Deep Healing: Ease the BTI and EM Wearout Crisis by Activating Recovery," In Proc. of 13th Workshop on Silicon Errors in Logic-System Effects (SELSE-13), Boston, MA, 2017. (Best Paper Award)
- [C11] X. Guo, M. Stan, "Enabling Wearout-Immune BEOL and FEOL with Active Rejuvenation," In Proc. of IEEE/ACM Workshop on Variability Modeling and Characterization (VMC), in conjunction with ICCAD, Austin, TX, 2016.
- [C12] X. Guo, M. Stan, "Work hard, sleep well Avoid irreversible IC wearout with proactive rejuvenation," In Proc. of the ACM/IEEE Asia and South Pacific Design Automation Conference (ASP-DAC), Macao, China, 2016.

- [C13] X. Guo, M. Stan, "MCPENS: Multiple-Critical-Path Embeddable NBTI Sensors for Dynamic Wearout Management," In Proc. of 11th Workshop on Silicon Errors in Logic-System Effects (SELSE-11), Austin, TX, 2015.
- [C14] X. Guo, W. Burleson, M. Stan, "Modeling and Experimental Demonstration of Accelerated Self-Healing Techniques," In *Proc. of ACM/IEEE Design Automation Conference* (DAC), San Fransisco, CA, 2014.
- [C15] Y. Zhao, Y. Yang, K. Mazumdar, X. Guo, M. Stan, "A Multi-Output on-Chip Switched-Capacitor DC-DC Converter for Near- and Subthreshold Power Modes," In Proc. of IEEE International Symposium on Circuits and Systems (ISCAS), Melbourne, Australia, 2014.

A.3 Posters

- [P1] D. Akella, X. Guo, M. Stan, B. H. Calhoun, "Enabling Post-Silicon Hold Time Closure by Tunable-Buffer Insertion," *Design Automation Conference* (DAC), Work-in-Progress (WIP) Poster Session, Austin, TX, 2017.
- [P2] X. Guo, M. Stan, "Exploring Aging Issues in Low-Power Internet of Things (IoT) Applications," *The 13th University of Virginia Engineering Research Symposium* (UVERS), Charlottesville, VA, 2017.
- [P3] X. Guo, M. Stan, "Towards Wearout-aware and Accelerated Self-healing Digital Systems," ACM Student Research Competition (SRC), in conjunction with ICCAD, Austin, TX, 2016.
- [P4] X. Guo, M. Stan, "CLASH Cross-Layer Accelerated Self-Healing: Circadian Rhythms for Resilient Electronic Systems," SRC System Level Design Review (SLD), Hillsboro, OR, 2016.
- [P5] X. Guo, M. Stan, "Towards Wearout-aware and Accelerated Self-Healing Digital Systems," ACM SIGDA Student Research Forum (SRF), in conjunction with Asia and South Pacific Design Automation Conference (ASP-DAC), Macao, China, 2016.
- [P6] A. Roelke, X. Guo, M. Stan, "Architecture Implications of Proactive Accelerated Wearout and Recovery," *Center for Future Architecture* (C-FAR) *Annual Review Poster Pitches*, Ann Arbor, MI, 2015.
- [P7] X. Guo, M. Stan, "Towards Wearout-aware and Accelerated Self-healing Digital Systems," ACM SIGDA PhD Forum, Design Automation Conference (DAC), San Fransisco, CA, 2015.
- [P8] M. Stan, X. Guo, A. Roelke, "Modeling and Experimental Demonstration of Accelerated Self-Healing Techniques in CMOS Circuits," *Government Microcircuit Applications & Critical Technology Conference* (GOMACTech), 2015.

- [P9] X. Guo, M. Stan, "Towards Aging-aware and Self-healing VLSI Chips and Systems," *The 11st University of Virginia Engineering Research Symposium* (UVERS), Charlottesville, VA, 2015.
- [P10] X. Guo, M. Stan, "Exploring Accelerated Self-Healing Techniques for Electronic Chips and Systems," 10th University of Virginia Engineering Research Symposium (UVERS), Charlottesville, VA, 2014.
- [P11] X. Guo, M. Stan, "Aging effects on FPGA chips and systems," A. Richard Newton Young Fellow Poster Session, Design Automation Conference (DAC), Austin, TX, 2013.

A.4 Presentations/Talks

- [T1] X. Guo, M. Stan, "Deep Healing: Ease the BTI and EM Wearout Crisis by Activating Recovery," In SRC TECHCON, Austin, TX, 2017. (Best-in-Session Award)
- [T2] M. El-Hadedy, X. Guo, W. Hwu, M. Stan, K. Skadron, "Crypt-Pi: A Light and Fast Crypto-Processor for IoT Applications," In SRC TECHCON, Austin, TX, 2017. (Bestin-Session Award)
- [T3] M. El-Hadedy, X. Guo, M. Stan, K. Skadron, W. Hwu, "R-NNPE: Reconfigurable Neural Network Processing Elements," In SRC TECHCON, Austin, TX, 2017.
- [T4] X. Guo, M. Stan, "Towards Wearout-Free Systems: A Self-Healing Strategy Enabled by Accelerated and Active Recovery," In SRC TECHCON, Austin, TX, 2016.
- [T5] M. El-Hadedy, X. Guo, M. Stan, K. Skadron, "Area-Efficient VLIW-based Programmable Processing Element for Crypto-Systems," In SRC TECHCON, Austin, TX, 2016.

Appendix B

Flow for Placing

Metastable-element-based BTI Sensors

In Chapter 4.4.3, we proposed a metastable-element-based BTI sensor for tracking both wearout and recovery. The sensor is small and digital-based, thus it can be embedded as an IP to be included in a standard synthesis flow. In this appendix, we will demonstrate the flow by instrumenting these sensors in a Johnson counter. This flow assumes the new scan cell which includes the sensors have been created by following the flow described in Chapter 4.4.3. Also this flow is exercised with a Synopsys flow, but similar methodology can be applied in Cadence environment.

B.1 New Top-down Design Flow with BTI Sensor Insertion

As shown in Figure B.1, the added steps for inserting BTI sensor IPs are during logic synthesis step. After the basic synthesis with Design Compiler, the DFT flow runs to replace the flip flops with scan cell. Following this, we need to update the netlist and replace these scan cell with the new scan cell which includes the sensor IP. During the physical design,



Fig. B.1 Updated Top-down Design Flow with BTI Sensor Insertion.

the new scancell IP needs to be added to the reference library so that the PnR tool is able to instantiate the design. The following section will demonstrate the flow step by step.

B.2 Demonstration in a Counter Design

As a proof-of-concept demonstration, we pick an 8-bit Johnson Counter as our target design for sensor insertion. The detailed steps are given in Figure B.2. After the regular DFT step, the placement strategy needs to be decided, it includes where to place the sensors, how many sensors are to be placed and what control signals need to be added. The updated netlist with the new scan cells is shown on the bottom half of the figure, where in this case, 3 out of 8 scan cells are replaced, the sensor inputs and control signals are also added in the signal list and are defined as inputs. After modifying the netlist from DFT, we continue P&R in IC Compiler (ICC) and the final layout with the sensor embedded is shown in Figure B.3.


Fig. B.2 Demonstration of the Sensor Insertion Flow in a Johnson Counter Design.

New Scan Cells with Sensor Embedded



Fig. B.3 The layout of the counter design after sensor insertion.

Appendix C

A Brief Overview of My Side Projects

The following projects are led by other individuals, but I contribute to some parts of them. They are briefly summarized here.

C.1 Post-silicon Hold Time Closure – Tunable Buffer Insertion

Hold-time is usually addressed in the design flow by inserting buffers in the datapath, and hold-induced failures are usually hard to fix after silicon is back. In this project, we propose a post-silicon hold time closure technique that uses tunable-buffers in the data-path instead of traditional-buffers. This enables post-silicon correction of hold violations and therefore, reduces the design effort in estimating design-time hold margins. We design a tunable buffer, demonstrate the tunable-buffer insertion strategy, and present a physical design flow using standard EDA tools. We verify this technique with measurements of a 130 nm test chip. A design-dependent hold slack improvement in the range of 103%-195% is achieved compared to the traditional buffering technique, with minimal power and area overhead. This technique also has the potential to reduce the number of buffers inserted for hold closure.

This work was in collaboration with Dr. Divya Akella Kamakshi (ECE department, University of Virginia), who also led this project. My contributions were: 1) Conducting the power analysis; 2) Helping with the timing closure; 3) System-level implementations. The work has been published in [C4] and [P1].

C.2 A 14nm Low-Vdd Heterogeneous RISC-V-based SoC

In this project, we design and implement in 14nm FinFET technology a heterogeneous RISC-V-based system, VELVET-bonus, capable of operating at very low voltage. VELVET-bonus encompasses an open-source RISC-V microprocessor (Rocket) with a tightly coupled machine learning accelerator (DANA), a power/resiliency management unit (PRIME), low voltage SRAMs, and instrumentation using new tools and languages via critical path monitors (CPMs) and power proxy sensors. The system is entirely RTL without the traditionally necessary architectural abstractions.

This project was in collaboration with my colleague from UVa and researchers from IBM T. J. Watson Research Center. My contributions in this tapeout are: 1) Backend implementation flow (RTL to GDS) evaluation and enhancement; 2) Top-level Integration; 3) Timing closure and signoff. The preliminary work has been presented in [C5] and [C9].

C.3 Programmable Processing Element for IoT Crypto Systems

This project introduced a novel programmable processing element (PPE) for various cryptographic systems that can be used in IoT applications. The design enables the programmability, thus supporting a wide range of bit-widths (such as 16, 32, and 64). It employs a very long instruction word (VLIW) architecture with an instruction set and mem-

ory hierarchy specialized for crypto-processing. Both FPGA and ASIC implementations demonstrate that the design utilizes a very tiny area and consumes very low power. For example, it takes only 227 slices for FPGA implementations to include 512-byte instruction and coefficient memory along with the computational unit by achieving a maximum clock frequency of 250MHz. For ASIC implementation (in 28/32nm technology), the design takes only $0.15mm^2$ of silicon area and consumes only 34.5W of total power while achieving a maximum frequency of 952MHz. To evaluate the effectiveness of the design in a larger system, we implement Blue Midnight Wish (BMW) hash function with the PPE. Compared to the previous BMW-512 implementation which stores the intermediate coefficients of the BMW-512 in 2048 bytes, the proposed design just uses 512 bytes. Meanwhile, we reduce the instruction memory size from 4864 bytes to 1792 bytes.

This project was led by Dr. Mohamed El-Hadedy (Cybersecurity Assistant Professor, ECE Department, California State Polytechnic University, Pomona). My contributions include: 1) ASIC implementation and evaluation of the design; 2) Area and power analysis. The work has been published in [C8].

C.4 Dual-Data Rate Transpose Memory

In this project, we propose a novel type of high-speed and area-efficient register-based transpose memory architecture enabled by reporting on both edges of the clock. The proposed new architecture, by using the double-edge triggered registers, doubles the throughput and increases the maximum frequency by avoiding some of the combinational circuit used in prior work. The proposed design is evaluated with both FPGA and ASIC flow in 28/32nm technology. The experimental results show that the proposed memory achieves almost $4 \times$ improvement in throughput while consuming 46% less area with the FPGA implementations compared to prior work. For ASIC implementations, it achieves more than 60% area reduction and at least $2 \times$ performance improvement while burning 60% less power compared

to other register-based designs implemented with the same flow. As an example, a proposed 8X8 transpose memory with 12-bit input/output resolution is able to achieve a throughput of 107.83*Gbps* at 647*MHz* by taking only 140 slices on a Virtex-7 Xilinx FPGA platform, and achieve a throughput of 88.2*Gbps* at 529*MHz* by taking $0.024mm^2$ silicon area for ASIC. The proposed transpose memory is integrated in both 2D-DCT and 2D-IDCT blocks for signal processing applications on the same FPGA platform. The new architecture allows a $3.5 \times$ speedup in performance for the 2D-DCT algorithm, compared to the previous work, while consuming 28% less area, and 2D-IDCT achieves a 3X speed-up while consuming 20% less area.

This project was led by Dr. Mohamed El-Hadedy (Cybersecurity Assistant Professor, ECE Department, California State Polytechnic University, Pomona). My contributions include: 1) ASIC evaluations of the design; 2) Scalability analysis; 3) Area and power analysis. The work has been published in [J3].

C.5 On-chip Power Regulation with Voltage Stacking

In this project, we present key experimental results of voltage stacking (VS) from a test setup using stacked loads (commercial-off-the-shelf FPGA chips as CMOS loads and passive resistors as resistive loads) and stacked Switched-Capacitor (SC) converters (dual-output push-pull SC converter, designed and fabricated in a 130*nm* bulk CMOS technology) to evaluate the performance benefits of this series-connected architectures. An architectural/circuit-level simulation framework has been created to run different benchmarks/applications in the SC converter assisted stacked loads architecture to demonstrate its ability to work with diverse nature of loads.

This project is led by Dr. Kaushik Mazumdar (ECE department, University of Virginia). My contributions are: 1) Helping with the testchip tapeout (DRC, LVS, ect.); 2) Setting up the FPGA test; 3) Helping with designing a multi-output on-chip switched-capacitor DC-DC Converter. The work has been published in [J5] and [C15].

Glossary

Subscripts

acce	Accelerated recovery
der	Derating
gs	Gate and source
life	Lifetime
nuc	Void Nucleation
osc	Oscillation
out	Output
OX	Oxide
rec	Recovery
ref	Reference
sat	Saturation
st	Stress
th	Threshold

Acronyms / Abbreviations

- ABB Adaptive Body Biasing
- ADE Analog Design Environment
- AR Accelerated Recovery
- AS Accelerated Stress
- ASIC Application Specific Integrated Circuit
- BEOL Back-end-of-line
- BOL Beginning of Lifetime
- BRAM Block Random-access Memory
- BTI Bias Temperature Instability
- CLASH Cross-Layer Accelerated Self-Healing
- CLM Channel Length Modulation
- CMOS Complementary Metal-oxide-semiconductor
- CUT Circuit Under Test
- Decap Decoupling Capacitor
- DFT Design for Test
- DIBL Drain-Induced-Barrier Lowering
- DRC Design Rule Check
- DVFS Dynamic Voltage Frequency Scaling

EDA	Electronic Design Automation
EDP	Energy Delay Product
EM	Electromigration
EOL	End of Lifetime
EOT	Equivalent Oxide Thickness
EPRO	M Erasable Programmable Read-only Memory
FBB	Forward Body Bias
FD-SC	DI Fully Depleted Silicon On Insulator
FEOL	Front-end-of-line
FinFE'	T Fin Field Effect Transistor
FO4	Fanout of 4
FPGA	Field-programmable Gate Array
GIDL	Gate-induced Drain Leakage
HCI	Hot Carrier Injection
HPLP	High-performance Low-power research group at the University of Virginia
I/O	Input and Output
IC	Integrated Circuit
IG	Independent Gate
IMP	Average Performance Improvement

IoE	Internet of Everything
IOE	Internet of Everything

- IoT Internet of Things
- IP Intellectual Properties
- IR-Drop Electrical potential difference between the two ends of a conducting phase during a current flow. This voltage drop across any resistance is the product of current (I) passing through resistance and resistance value (R).
- IR Irreversible Wearout
- ISA Instruction Set Architecture
- LELE DP Litho-Etch-Litho-Etch Double Patterning
- LER line-edge Roughness
- LUT Look-up-table
- LV Low Voltage
- MCU Micro-controller Unit
- MEOL Middle-end-of-line
- MRAM Magnetic Random Access Memory
- MTJ Magnetic Tunnel Junction
- MTTF Mean Time To Failure
- MUX Multiplexer
- NBTI Negative-bias Temperature Instability
- NMOS N-channel Metal–Oxide–Semiconductor Field-effect Transistor

N TT 7N	N T 1		<i>r</i> •
N N N N	Non vo	lotila N	amoriac.
	1 NOII - VOI		
	1,011,00		101101100

- PBTI Positive-bias Temperature Instability
- PC Personal Computer
- PCRAM Phase-change Random Access Memory
- PDC Pulsed DC
- PDK Process Design Kit
- PDN Power Delivery Network
- PLB Programmable Logic Block
- PMOS P-channel Metal–Oxide–Semiconductor Field-effect Transistor
- PMU Power Management Unit
- PnR Place and Route
- POI Path of Interest
- PPA Power Performance Area
- PUN Pull Up Network
- **RDF** Random Dopant Fluctuations
- RD Reaction-Diffusion
- rob Reorder Buffer
- RO Ring Oscillator
- RRAM Resistive Random Access Memory

RTL	Register-transfer Level	
SADP	Self-aligned Double Patterning	
SC	Switched-Capacitor	
SEM	Scanning Electron Microscope	
SG	Shorted Gate	
SID SADP Spacer-is-Dielectric Double Patterning		
SIM S	ADP Spacer-is-Metal Self-Aligned Double Patterning	
SoC	System-on-Chip	
SOI	Silicon On Insulator	
SPEC	Specifications	
SPICE Simulation Program with Integrated Circuit Emphasis		
SPI	Serial Peripheral Interface	
SRAM	I Static Random-access Memory	
STI	Shallow Trench Isolation	
ST	Sleep Transistor	
STT	Spin-transfer Torque	
TCAD	Technology Computer-Aided Design	
TD	Trapping-Detrapping	
TEI	Thermal Effect Inversion	

TSV	Through Silicon Vias
TTF	Time To Failure
TT	Typical-Typical
ULP	Ultra-low Power
UV	Ultraviolet
VLSI	Very-large-scale Integration
VS	Voltage Stacking
VTC	Voltage Transfer Curve

Bibliography

- [1] Semiconductor Engineering Reliability Challenges In 16nm FinFET Design:. http://semiengineering.com/reliability-challenges-16nm-finfet-design/.
- [2] Jaume Abella and Xavier Vera. Electromigration for microarchitects. *ACM Computing Surveys (CSUR)*, 42(2):9, 2010.
- [3] Jaume Abella, Xavier Vera, and Antonio Gonzalez. Penelope: The nbti-aware processor. In *Microarchitecture*, 2007. *MICRO* 2007. 40th Annual IEEE/ACM International Symposium on, pages 85–96. IEEE, 2007.
- [4] Jaume Abella, Xavier Vera, Osman S Unsal, Oguz Ergin, Antonio González, and James W Tschanz. Refueling: Preventing wire degradation due to electromigration. *IEEE micro*, (6):37–46, 2008.
- [5] Thomas Aichinger, Michael Nelhiebel, and Tibor Grasser. On the temperature dependence of nbti recovery. *Microelectronics Reliability*, 48(8):1178–1184, 2008.
- [6] Rob Aitken, Ethan H Cannon, Mondira Pant, and Mehdi B Tahoori. Resiliency challenges in sub-10nm technologies. In VLSI Test Symposium (VTS), 2015 IEEE 33rd, pages 1–4. IEEE, 2015.
- [7] Robert Aitken, Greg Yeric, Brian Cline, Saurabh Sinha, Lucian Shifren, Imran Iqbal, and Vikas Chandra. Physical Design and FinFETs. In *Proceedings of the 2014 on International symposium on physical design*, pages 65–68. ACM, 2014.
- [8] Hossein Karimiyan Alidash, Andrea Calimera, Alberto Macii, Enrico Macii, and Massimo Poncino. On-chip nbti and pbti tracking through an all-digital aging monitor architecture. In *PATMOS*, pages 155–165. Springer, 2012.
- [9] Massimo Alioto. Enabling the Internet of Things: From Integrated Circuits to Integrated Systems. Springer, 2017.
- [10] A. Amouri, J. Hepp, and M. Tahoori. Built-in self-heating thermal testing of fpgas. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, PP (99):1–1, 2016. ISSN 0278-0070. doi: 10.1109/TCAD.2015.2512905.

- [11] Abdulazim Amouri, Jochen Hepp, and Mehdi Tahoori. Self-heating thermal-aware testing of fpgas. In VLSI Test Symposium (VTS), 2014 IEEE 32nd, pages 1–6. IEEE, 2014.
- [12] Digilent Analog Discovery 100MS/s USB Oscilloscope & Logic Analyzer:. http://store.digilentinc.com/ analog-discovery-100msps-usb-oscilloscope-logic-analyzer-limited-time/.
- [13] Dean Michael Ancajas, Koushik Chakraborty, and Sanghamitra Roy. Proactive aging management in heterogeneous nocs through a criticality-driven routing approach. In *Proceedings of the Conference on Design, Automation and Test in Europe*, pages 1032–1037. EDA Consortium, 2013.
- [14] D Angot, V Huard, X Federspiel, F Cacho, and A Bravaix. Bias temperature instability and hot carrier circuit ageing simulations specificities in UTBB FDSOI 28nm node. In *Reliability Physics Symposium (IRPS), 2013 IEEE International*, pages 5D–2. IEEE, 2013.
- [15] Rizwan A Ashraf, Navid Khoshavi, Ahmad Alzahrani, Ronald F DeMara, Saman Kiamehr, and Mehdi B Tahoori. Area-energy tradeoffs of logic wear-leveling for bti-induced aging. In *Proceedings of the ACM International Conference on Computing Frontiers*, pages 37–44. ACM, 2016.
- [16] Aditya Bansal and Jae-Joon Kim. Power napping technique for accelerated negative bias temperature instability (nbti) and/or positive bias temperature instability (pbti) recovery, July 21 2015. US Patent 9086865.
- [17] Emanuele Baravelli, Malgorzata Jurczak, Nicolò Speciale, Kristin De Meyer, and Abhisek Dixit. Impact of LER and Random Dopant Fluctuations on FinFET Matching Performance. *IEEE transactions on nanotechnology*, 7(3):291–298, 2008.
- [18] Robert Baumann. Soft errors in advanced computer systems. *IEEE Design & Test of Computers*, 22(3):258–266, 2005.
- [19] A Benabdelmoumene, B Djezzar, A Chenouf, H Tahi, B Zatout, and M Kechouane. On the turn-around phenomenon in n-mos transistors under nbti conditions. *Solid-State Electronics*, 121:34–40, 2016.
- [20] Kerry Bernstein, David J Frank, Anne E Gattiker, Wilfried Haensch, Brian L Ji, Sani R Nassif, Edward J Nowak, Dale J Pearson, and Norman J Rohrer. High-performance cmos variability in the 65-nm regime and beyond. *IBM journal of research and development*, 50(4.5):433–449, 2006.
- [21] Debajit Bhattacharya and Niraj K Jha. FinFETs: From Devices to Architectures. *Advances in Electronics*, 2014, 2014.
- [22] Andy Biddle and Jason ST Chen. FinFET Technology-Understanding and Productizing a New Transistor. *A joint whitepaper from TSMC and Synopsys*, 2013.

- [23] Christian Bienia. Benchmarking modern multiprocessors. Princeton University, 2011.
- [24] Nathan Binkert, Bradford Beckmann, Gabriel Black, Steven K Reinhardt, Ali Saidi, Arkaprava Basu, Joel Hestness, Derek R Hower, Tushar Krishna, Somayeh Sardashti, et al. The gem5 simulator. ACM SIGARCH Computer Architecture News, 39(2):1–7, 2011.
- [25] James R Black. Electromigration—a brief survey and some recent results. *IEEE Transactions on Electron Devices*, 16(4):338–347, 1969.
- [26] Illan A Blech. Electromigration in thin aluminum films on titanium nitride. *Journal* of Applied Physics, 47(4):1203–1208, 1976.
- [27] Paul Bogdan, Siddharth Garg, and Umit Y Ogras. Energy-efficient computing from systems-on-chip to micro-server and data centers. In *Green Computing Conference* and Sustainable Computing Conference (IGSC), 2015 Sixth International, pages 1–6. IEEE, 2015.
- [28] Pradip Bose, Jeonghee Shin, and Victor Zyuban. Method for extending lifetime reliability of digital logic devices through removal of aging mechanisms, February 10 2009. US Patent 7,489,161.
- [29] Pradip Bose, Jeonghee Shin, and Victor Zyuban. Method for extending lifetime reliability of digital logic devices through reversal of aging mechanisms, February 3 2009. US Patent 7,486,107.
- [30] D. Burnett, S. Parihar, H. Ramamurthy, and S. Balasubramanian. Finfet sram design challenges. In 2014 IEEE International Conference on IC Design Technology, pages 1–4, May 2014. doi: 10.1109/ICICDT.2014.6838606.
- [31] Ermao Cai and Diana Marculescu. TEI-turbo: Temperature Effect Inversion-aware Turbo Boost for FinFET-based Multi-core Systems. In *Proceedings of the IEEE/ACM International Conference on Computer-Aided Design*, pages 500–507. IEEE Press, 2015.
- [32] Ermao Cai, Dimitrios Stamoulis, and Diana Marculescu. Exploring Aging Deceleration in FinFET-based Multi-core Systems. In *Computer-Aided Design (ICCAD)*, 2016 IEEE/ACM International Conference on, pages 1–8. IEEE, 2016.
- [33] Benton H Calhoun, Yu Cao, Xin Li, Ken Mai, Lawrence T Pileggi, Rob A Rutenbar, and Kenneth L Shepard. Digital Circuit Design Challenges and Opportunities in the Era of Nanoscale CMOS. *Proceedings of the IEEE*, 96(2):343–365, 2008.
- [34] Andrea Calimera, Enrico Macii, and Massimo Poncino. Nbti-aware sleep transistor design for reliable power-gating. In *Proceedings of the 19th ACM Great Lakes* symposium on VLSI, pages 333–338. ACM, 2009.

- [35] Andrea Calimera, Enrico Macii, and Massimo Poncino. Design techniques for nbtitolerant power-gating architectures. *Circuits and Systems II: Express Briefs, IEEE Transactions on*, 59(4):249–253, 2012.
- [36] Andrea Calimera, Alberto Macii, Enrico Macii, and Massimo Poncino. Power-gating for leakage control and beyond. In *Circuit Design for Reliability*, pages 175–205. Springer, 2015.
- [37] Yu Cao, Jyothi Velamala, Ketul Sutaria, Mike Shuo-Wei Chen, Jonathan Ahlbin, Ivan Sanchez Esqueda, Michael Bajura, and Michael Fritze. Cross-layer modeling and simulation of circuit reliability. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 33(1):8–23, 2014.
- [38] Nicholas P Carter, Helia Naeimi, and Donald S Gardner. Design techniques for crosslayer resilience. In *Proceedings of the Conference on Design, Automation and Test in Europe*, pages 1023–1028. European Design and Automation Association, 2010.
- [39] Tuck-Boon Chan, John Sartori, Puneet Gupta, and Rakesh Kumar. On the efficacy of nbti mitigation techniques. In *Design, Automation & Test in Europe Conference & Exhibition (DATE), 2011*, pages 1–6. IEEE, 2011.
- [40] Wei-Ting Chan, Andrew B Kahng, and Siddhartha Nath. Methodology for electromigration signoff in the presence of adaptive voltage scaling. In System Level Interconnect Prediction (SLIP), 2014 ACM/IEEE International Workshop on, pages 1–7. IEEE, 2014.
- [41] Wen-Teng Chang, Shih-Wei Lin, Cheng-Ting Shih, and Wen-Kuan Yeh. Back Bias Modulation of UTBB FDSOI, Bulk FinFET, and SOI FinFET. In *Nanoelectronics Conference (INEC), 2016 IEEE International*, pages 1–2. IEEE, 2016.
- [42] GCKY Chen, KY Chuah, MF Li, Daniel SH Chan, CH Ang, JZ Zheng, Y Jin, and DL Kwong. Dynamic nbti of pmos transistors and its impact on device lifetime. In *Reliability Physics Symposium Proceedings*, 2003. 41st Annual. 2003 IEEE International, pages 196–202. IEEE, 2003.
- [43] Xiaoming Chen, Yu Wang, Yu Cao, Yuchun Ma, and Huazhong Yang. Variation-aware supply voltage assignment for simultaneous power and aging optimization. Very Large Scale Integration (VLSI) Systems, IEEE Transactions on, 20(11):2143–2147, 2012.
- [44] E. Cheng, J. Abraham, P. Bose, A. Buyuktosunoglu, K. Campbell, D. Chen, C. Y. Cher, H. Cho, B. Le, K. Lilja, S. Mirkhani, K. Skadron, M. Stan, L. Szafaryn, C. Vezyrtzis, and S. Mitra. Cross-layer resilience in low-voltage digital systems: Key insights. In 2017 IEEE International Conference on Computer Design (ICCD), pages 593–596, Nov 2017.
- [45] Thomas Chiarella, Liesbeth Witters, Abdelkarim Mercha, Christoph Kerner, Michal Rakowski, Claude Ortolland, L-Å Ragnarsson, Bertrand Parvais, Ari De Keersgieter,

Stefan Kubicek, et al. Benchmarking SOI and Bulk FinFET Alternatives for PLANAR CMOS Scaling Succession. *Solid-State Electronics*, 54(9):855–860, 2010.

- [46] Minki Cho, Stephen T Kim, Carlos Tokunaga, Charles Augustine, Jaydeep P Kulkarni, Krishnan Ravichandran, James W Tschanz, Muhammad M Khellah, and Vivek De. Postsilicon voltage guard-band reduction in a 22 nm graphics execution core using adaptive voltage scaling and dynamic power gating. *IEEE Journal of Solid-State Circuits*, 52(1):50–63, 2017.
- [47] Seongjae Cho, Jung Hoon Lee, Shinichi O'uchi, Kazuhiko Endo, Meishoku Masahara, and Byung-Gook Park. Design of SOI FinFET on 32nm Technology Node for Low Standby Power (LSTP) Operation Considering Gate-induced Drain Leakage (GIDL). *Solid-State Electronics*, 54(10):1060–1065, 2010.
- [48] Lawrence T Clark, Vinay Vashishtha, Lucian Shifren, Aditya Gujja, Saurabh Sinha, Brian Cline, Chandarasekaran Ramamurthy, and Greg Yeric. ASAP7: A 7-nm FinFET Predictive Process Design Kit. *Microelectronics Journal*, 53:105–115, 2016.
- [49] Middle East Technical University Computer Simulation Laboratory (CSL):. http: //www.csl.mete.metu.edu.tr/Electromigration/emig.htm.
- [50] Maxim Switched-Capacitor Voltage Converters MAX1044/ICL7660 datasheet:. https://www.maximintegrated.com/en/products/power/charge-pumps/ICL7660.html.
- [51] M Denais, C Parthasarathy, G Ribes, Y Rey-Tauriac, N Revil, A Bravaix, V Huard, and F Perrier. On-the-fly characterization of nbti in ultra-thin gate oxide pmosfet's. In *Electron Devices Meeting, 2004. IEDM Technical Digest. IEEE International*, pages 109–112. IEEE, 2004.
- [52] Boualem Djezzar, Hakim Tahi, Abdelmadjid Benabdelmoumene, Amel Chenouf, Mohamed Goudjil, and Youcef Kribes. On the permanent component profiling of the negative bias temperature instability in p-mosfet devices. *Solid-State Electronics*, 106: 54–62, 2015.
- [53] J El Husseini, A Subirats, X Garros, A Makoseij, O Thomas, G Reimbold, V Huard, F Cacho, and X Federspiel. Accurate modeling of dynamic variability of sram cell in 28 nm fdsoi technology. In *Microelectronic Test Structures (ICMTS), 2014 International Conference on*, pages 41–46. IEEE, 2014.
- [54] Semiconductor Engineering. Chip Aging Accelerates. http://semiengineering.com/ chip-aging-accelerates/#.WoTxbxFVVp0.email, 2018.
- [55] Hadi Esmaeilzadeh, Emily Blem, Renee St Amant, Karthikeyan Sankaralingam, and Doug Burger. Dark silicon and the end of multicore scaling. In *Computer Architecture* (ISCA), 2011 38th Annual International Symposium on, pages 365–376. IEEE, 2011.

- [56] Gregory G Faust, Runjie Zhang, Kevin Skadron, Mircea R Stan, and Brett H Meyer. ArchFP: Rapid prototyping of pre-RTL floorplans. In VLSI and System-on-Chip (VLSI-SoC), 2012 IEEE/IFIP 20th International Conference on, pages 183–188. IEEE, 2012.
- [57] Raoul Fernandez, Ben Kaczer, Axel Nackaerts, Steven Demuynck, R Rodriguez, Montserat Nafria, and Guido Groeseneken. Ac nbti studied in the 1 hz–2 ghz range on dedicated on-chip cmos circuits. In *Electron Devices Meeting*, 2006. *IEDM'06*. *International*, pages 1–4. IEEE, 2006.
- [58] FinFET, Multi-Patterning Aware Place, and Route Implementation:. http://go.mentor. com/4h_c2.
- [59] Re-Engineering The FinFET:. http://semiengineering.com/re-engineering-the-finfet/.
- [60] Farshad Firouzi, Fangming Ye, Krishnendu Chakrabarty, and Mehdi B Tahoori. Agingand variation-aware delay monitoring using representative critical path selection. *ACM Transactions on Design Automation of Electronic Systems (TODAES)*, 20(3):39, 2015.
- [61] Jacopo Franco, Salvatore Graziano, Ben Kaczer, Felice Crupi, L-Å Ragnarsson, Tibor Grasser, and Guido Groeseneken. Bti reliability of ultra-thin eot mosfets for subthreshold logic. *Microelectronics Reliability*, 52(9):1932–1935, 2012.
- [62] Hasse Fredriksson and Ulla Akerlind. *Physics of functional materials*. John Wiley & Sons, 2008.
- [63] FreePDK45 from NCSU:. https://www.eda.ncsu.edu/wiki/FreePDK45:Contents.
- [64] Anshul Gandhi, Mor Harchol-Balter, and Michael A Kozuch. Are sleep states effective in data centers? In *Green Computing Conference (IGCC), 2012 International*, pages 1–10. IEEE, 2012.
- [65] Brad D Gaynor and Soha Hassoun. Fin Shape Impact on FinFET Leakage with Application to Multithreshold and Ultralow-leakage FinFET design. *IEEE Transactions on Electron Devices*, 61(8):2738–2744, 2014.
- [66] Bradley Geden. Understand and avoid electromigration (em) & ir-drop in custom ip blocks. *Synopsys White Paper*, pages 1–6, 2011.
- [67] Mohammad Saber Golanbari, Nour Sayed, Mojtaba Ebrahimi, Mohammad Hadi Moshrefpour Esfahany, Saman Kiamehr, and Mehdi B Tahoori. Aging-aware coding scheme for memory arrays. In *Test Symposium (ETS), 2017 22nd IEEE*, pages 1–6. IEEE, 2017.
- [68] T Grasser, M Waltl, G Rzepa, W Goes, Y Wimmer, A-M El-Sayed, AL Shluger, H Reisinger, and B Kaczer. The "permanent" component of nbti revisited: Saturation, degradation-reversal, and annealing. In *Reliability Physics Symposium (IRPS), 2016 IEEE International*, pages 5A–2. IEEE, 2016.

- [69] Tibor Grasser, B Kaczer, W Goes, Th Aichinger, Ph Hehenberger, and M Nelhiebel. A two-stage model for negative bias temperature instability. In *Reliability Physics Symposium, 2009 IEEE International*, pages 33–44. IEEE, 2009.
- [70] Tibor Grasser, Th Aichinger, Gregor Pobegen, Hans Reisinger, P-J Wagner, Jacopo Franco, M Nelhiebel, and Ben Kaczer. The 'permanent' component of nbti: composition and annealing. In *Reliability Physics Symposium (IRPS), 2011 IEEE International*, pages 6A–2. IEEE, 2011.
- [71] Jie Gu, John Keane, Sachin Sapatnekar, and Chris Kim. Width Quantization Aware FinFET Circuit Design. In *Custom Integrated Circuits Conference*, 2006. CICC'06. IEEE, pages 337–340. IEEE, 2006.
- [72] Jayavardhana Gubbi, Rajkumar Buyya, Slaven Marusic, and Marimuthu Palaniswami. Internet of Things (IoT): A vision, architectural elements, and future directions. *Future generation computer systems*, 29(7):1645–1660, 2013.
- [73] Saket Gupta and Sachin S Sapatnekar. Gnomo: Greater-than-nominal v dd operation for bti mitigation. In *Design Automation Conference (ASP-DAC), 2012 17th Asia and South Pacific*, pages 271–276. IEEE, 2012.
- [74] Saket Gupta and Sachin S Sapatnekar. Employing circadian rhythms to enhance power and reliability. ACM Transactions on Design Automation of Electronic Systems (TODAES), 18(3):38, 2013.
- [75] Christine S Hau-Riege. An introduction to cu electromigration. *Microelectronics Reliability*, 44(2):195–205, 2004.
- [76] Kai He. *Parallel CAD Algorithms and Hardware Security for VLSI Systems*. PhD thesis, University of California, Riverside, 2016.
- [77] Kai He, Xin Huang, and Sheldon X-D Tan. Em-based on-chip aging sensor for detection and prevention of counterfeit and recycled ics. In *Computer-Aided Design* (ICCAD), 2015 IEEE/ACM International Conference on, pages 146–151. IEEE, 2015.
- [78] Jörg Henkel, Lars Bauer, Nikil Dutt, Puneet Gupta, Sani Nassif, Muhammad Shafique, Mehdi Tahoori, and Norbert Wehn. Reliable on-chip systems in the nano-era: lessons learnt and future trends. In *Proceedings of the 50th Annual Design Automation Conference*, page 99. ACM, 2013.
- [79] Jorg Henkel, Heba Khdr, Santiago Pagani, and Muhammad Shafique. New trends in dark silicon. In *Design Automation Conference (DAC)*, 2015 52nd ACM/EDAC/IEEE, pages 1–6. IEEE, 2015.
- [80] Hyejeong Hong, Jaeil Lim, Hyunyul Lim, and Sungho Kang. Lifetime reliability enhancement of microprocessors: Mitigating the impact of negative bias temperature instability. *ACM Computing Surveys (CSUR)*, 48(1):9, 2015.

- [81] C Calvin Hu. Modern Semiconductor Devices for Integrated Circuits. *Part 7: MOS-FETs in ICs Scaling, Leakage, and Other Topics*, 2011.
- [82] Lin Huang and Qiang Xu. Characterizing the lifetime reliability of manycore processors with core-level redundancy. In *Computer-Aided Design (ICCAD)*, 2010 IEEE/ACM International Conference on, pages 680–685. IEEE, 2010.
- [83] Wei Huang, Shougata Ghosh, Sivakumar Velusamy, Karthik Sankaranarayanan, Kevin Skadron, and Mircea R Stan. HotSpot: A compact thermal modeling methodology for early-stage VLSI design. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 14(5):501–513, 2006.
- [84] Xin Huang, Tan Yu, Valeriy Sukharev, and Sheldon X-D Tan. Physics-based electromigration assessment for power grid networks. In *Proceedings of the 51st Annual Design Automation Conference*, pages 1–6. ACM, 2014.
- [85] Xin Huang, Valeriy Sukharev, Taeyoung Kim, Haibao Chen, and Sheldon X-D Tan. Electromigration recovery modeling and analysis under time-dependent current and temperature stressing. In 2016 21st Asia and South Pacific Design Automation Conference (ASP-DAC), pages 244–249. IEEE, 2016.
- [86] Xin Huang, Valeriy Sukharev, Taeyoung Kim, and Sheldon X-D Tan. Dynamic electromigration modeling for transient stress evolution and recovery under timedependent current and temperature stressing. *Integration, the VLSI Journal*, 2016.
- [87] V Huard, M Denais, and C Parthasarathy. Nbti degradation: From physical mechanisms to modelling. *Microelectronics Reliability*, 46(1):1–23, 2006.
- [88] Vincent Huard, Florian Cacho, Xavier Federspiel, and Pascal Mora. Hot-carrier injection degradation in advanced cmos nodes: a bottom-up approach to circuit and system reliability. In *Hot Carrier Degradation in Semiconductor Devices*, pages 401–444. Springer, 2015.
- [89] Lattice Semiconductor iCE40 HX-Series Ultra Low-Power mobile FPGA Family Datasheet:. http://www.latticesemi.com/Products/FPGAandCPLD/iCE40.aspx.
- [90] Mitsuhiko Igarashi, Kan Takeuchi, Takeshi Okagaki, Koji Shibutani, Hiroaki Matsushita, and Koji Nii. An on-die digital aging monitor against hei and xbti in 16 nm fin-fet bulk emos technology. In *European Solid-State Circuits Conference (ESSCIRC)*, *ESSCIRC 2015-41st*, pages 112–115. IEEE, 2015.
- [91] Gopal Singh Jamnal, Xiaodong Liu, Lu Fan, and Muthu Ramachandran. Cognitive Internet of Everything (CIoE): State of the Art and Approaches. In *Emerging Trends and Applications of the Internet of Things*, pages 277–309. IGI Global, 2017.
- [92] Hai Jiang, SangHoon Shin, Xiaoyan Liu, Xing Zhang, and Muhammad Ashraful Alam. The Impact of Self-Heating on HCI Reliability in High-Performance Digital Circuits. *IEEE Electron Device Letters*, 38(4):430–433, 2017.

- [93] Norihiro Kamae, Akira Tsuchiya, and Hidetoshi Onodera. A body bias generator compatible with cell-based design flow for within-die variability compensation. In *Solid State Circuits Conference (A-SSCC), 2012 IEEE Asian*, pages 389–392. IEEE, 2012.
- [94] Norihiro Kamae, AKM Mahfuzul Islam, Akira Tsuchiya, and Hidetoshi Onodera. A body bias generator with wide supply-range down to threshold voltage for within-die variability compensation. In *Solid-State Circuits Conference (A-SSCC)*, 2014 IEEE Asian, pages 53–56. IEEE, 2014.
- [95] Byungseok Kang, Daecheon Kim, and Hyunseung Choo. Internet of Everything: A large-scale autonomic IoT gateway. *IEEE Transactions on Multi-Scale Computing Systems*, 2017.
- [96] Kunhyuk Kang, Keejong Kim, Ahmad E Islam, Muhammad Alam, Kaushik Roy, et al. Characterization and estimation of circuit reliability degradation under nbti using on-line iddq measurement. In *Design Automation Conference*, 2007. DAC'07. 44th ACM/IEEE, pages 358–363. IEEE, 2007.
- [97] Kunhyuk Kang, Saakshi Gangwal, Sang Phill Park, and Kaushik Roy. NBTI induced performance degradation in logic and memory circuits: how effectively can we approach a reliability solution? In *Proceedings of the 2008 Asia and South Pacific Design Automation Conference*, pages 726–731. IEEE Computer Society Press, 2008.
- [98] Eric Karl, Dennis Sylvester, and David Blaauw. Analysis of system-level reliability factors and implications on real-time monitoring methods for oxide breakdown device failures. In *Quality Electronic Design*, 2008. ISQED 2008. 9th International Symposium on, pages 391–395. IEEE, 2008.
- [99] Anastasios A Katsetos. Negative bias temperature instability (nbti) recovery with bake. *Microelectronics Reliability*, 48(10):1655–1659, 2008.
- [100] Xrysovalantis Kavousianos, Krishnendu Chakrabarty, Arvind Jain, and Rubin Parekhji. Test schedule optimization for multicore socs: Handling dynamic voltage scaling and multiple voltage islands. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 31(11):1754–1766, 2012.
- [101] Jamil Kawa. Designing with FinFETs: The Opportunities and the Challenges. In *Synopsys White Paper*, pages 1–8. Synopsys, 2012.
- [102] John Keane, Tae-Hyoung Kim, and Chris H Kim. An on-chip nbti sensor for measuring pmos threshold voltage degradation. *IEEE transactions on very large scale integration* (VLSI) systems, 18(6):947–956, 2010.
- [103] John Keane, Xiaofei Wang, Devin Persaud, and Chris H Kim. An all-in-one silicon odometer for separately monitoring hci, bti, and tddb. *IEEE Journal of Solid-State Circuits*, 45(4):817–829, 2010.

- [104] Pranita Kerber, Qintao Zhang, Siyuranga Koswatta, and Andres Bryant. GIDL in Doped and Undoped FinFET Devices for Low-leakage Applications. *IEEE Electron Device Letters*, 34(1):6–8, 2013.
- [105] Navid Khoshavi, Rizwan A Ashraf, and Ronald F DeMara. Applicability of powergating strategies for aging mitigation of CMOS logic paths. In *Circuits and Systems* (MWSCAS), 2014 IEEE 57th International Midwest Symposium on, pages 929–932. IEEE, 2014.
- [106] Navid Khoshavi, Rizwan A Ashraf, Ronald F DeMara, Saman Kiamehr, Fabian Oboril, and Mehdi B Tahoori. Contemporary CMOS aging mitigation techniques: Survey, taxonomy, and methods. *Integration, the VLSI Journal*, 59:10–22, 2017.
- [107] Saman Kiamehr, Farshad Firouzi, Mojtaba Ebrahimi, and Mehdi B Tahoori. Agingaware standard cell library design. In *Proceedings of the conference on Design*, *Automation & Test in Europe*, page 261. European Design and Automation Association, 2014.
- [108] Tae-Hyoung Kim, Randy Persaud, and Chris H Kim. Silicon odometer: An on-chip reliability monitor for measuring frequency degradation of digital circuits. *Solid-State Circuits, IEEE Journal of*, 43(4):874–880, 2008.
- [109] Taeyoung Kim, Xin Huang, Hai-Bao Chen, Valeriy Sukharev, and Sheldon X-D Tan. Learning-based dynamic reliability management for dark silicon processor considering em effects. In *Design, Automation & Test in Europe Conference & Exhibition (DATE),* 2016, pages 463–468. IEEE, 2016.
- [110] Tony Tae-Hyoung Kim, Pong-Fei Lu, Keith A Jenkins, and Chris H Kim. A ringoscillator-based reliability monitor for isolated measurement of nbti and pbti in highk/metal gate technology. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 23(7):1360–1364, 2015.
- [111] Woongrae Kim, Taizhi Liu, and Linda Milor. On-line monitoring of system health using on-chip srams as a wearout sensor. In On-Line Testing and Robust System Design (IOLTS), 2017 IEEE 23rd International Symposium on, pages 253–258. IEEE, 2017.
- [112] Giray Kömürcü, Ali Emre Pusane, and Günhan Dündar. Effects of aging and compensation mechanisms in ordering based ro-pufs. *Integration, the VLSI Journal*, 52: 71–76, 2016.
- [113] Abhishek Koneru, Arunkumar Vijayan, Krishnendu Chakrabarty, and Mehdi B Tahoori. Fine-grained aging prediction based on the monitoring of run-time stress using dft infrastructure. In *Computer-Aided Design (ICCAD), 2015 IEEE/ACM International Conference on*, pages 51–58. IEEE, 2015.

- [114] Sanjay V Kumar, Chris H Kim, and Sachin S Sapatnekar. Adaptive techniques for overcoming performance degradation due to aging in digital circuits. In *Proceedings* of the Asia and South Pacific Design Automation Conference, pages 284–289. IEEE Press, 2009.
- [115] Mark LaPedus. Interconnect challenges rising: Resistance and capacitance drive need for new materials and approaches. <u>http://semiengineering.com/</u> interconnect-challenges-grow-3/, 2016.
- [116] Chi-Shuen Lee, Brian Cline, Saurabh Sinha, Greg Yeric, and H-S Philip Wong. 32-bit Processor Core at 5-nm Technology: Analysis of Transistor and Interconnect Impact on VLSI System Performance. In *Electron Devices Meeting (IEDM), 2016 IEEE International*, pages 28–3. IEEE, 2016.
- [117] Jong-Ho Lee. Bulk FinFETs: Design at 14 nm Node and Key Characteristics. In Nano Devices and Circuit Techniques for Low-Energy Applications and Energy Harvesting, pages 33–64. Springer, 2016.
- [118] Ki-Don Lee. Electromigration recovery and short lead effect under bipolar-and unipolar-pulse current. In *Reliability Physics Symposium (IRPS), 2012 IEEE International*, pages 6B–3. IEEE, 2012.
- [119] Ming-Chao Lee, Yu-Guang Chen, Ding-Kei Huang, and Shih-Chieh Chang. Nbtiaware power gating design. In *Design Automation Conference (ASP-DAC), 2011 16th Asia and South Pacific*, pages 609–614. IEEE, 2011.
- [120] Tsung-Lin Lee, Chih Chieh Yeh, Chang-Yun Chang, and Feng Yuan. FinFETs With Different Fin Heights, August 23 2016. US Patent 9,425,102.
- [121] Woojoo Lee, Yanzhi Wang, Tiansong Cui, Shahin Nazarian, and Massoud Pedram. Dynamic Thermal Management for FinFET-based Circuits Exploiting the Temperature Effect Inversion Phenomenon. In *Proceedings of the 2014 international symposium* on Low power electronics and design, pages 105–110. ACM, 2014.
- [122] Scott Lerner and Baris Taskin. Workload-aware ASIC flow for lifetime improvement of multi-core IoT processors. In *Quality Electronic Design (ISQED)*, 2017 18th International Symposium on, pages 379–384. IEEE, 2017.
- [123] Jiangyi Li and Mingoo Seok. Robust and in-situ self-testing technique for monitoring device aging effects in pipeline circuits. In *Proceedings of the 51st Annual Design Automation Conference*, pages 1–6. ACM, 2014.
- [124] Lin Li, Youtao Zhang, Jun Yang, and Jianhua Zhao. Proactive nbti mitigation for busy functional units in out-of-order microprocessors. In *Proceedings of the Conference* on Design, Automation and Test in Europe, pages 411–416. European Design and Automation Association, 2010.

- [125] Sheng Li, Jung Ho Ahn, Richard D Strong, Jay B Brockman, Dean M Tullsen, and Norman P Jouppi. McPAT: an integrated power, area, and timing modeling framework for multicore and manycore architectures. In *Microarchitecture, 2009. MICRO-42*. *42nd Annual IEEE/ACM International Symposium on*, pages 469–480. IEEE, 2009.
- [126] Jens Lienig and Göran Jerke. Electromigration-aware physical design of integrated circuits. In VLSI Design, 2005. 18th International Conference on, pages 77–82. IEEE, 2005.
- [127] Chi-Hung Lin, Chia-Shiang Chen, Yu-He Chang, Yu-Ting Zhang, Shang-Rong Fang, Santanu Santra, and Rung-Bin Lin. Design Space Exploration of FinFETs with Double Fin Heights for Standard Cell Library. In VLSI (ISVLSI), 2016 IEEE Computer Society Annual Symposium on, pages 673–678. IEEE, 2016.
- [128] MH Lin and AS Oates. Ac and pulsed-dc stress electromigration failure mechanisms in cu interconnects. In *Interconnect Technology Conference (IITC)*, 2013 IEEE International, pages 1–3. IEEE, 2013.
- [129] Yongxun Liu, Kenichi Ishii, Meishoku Masahara, Toshiyuki Tsutsumi, Hidenori Takashima, Hiromi Yamauchi, and Eiichi Suzuki. Cross-sectional Channel Shape Dependence of Short-channel Effects in Fin-type Double-gate Metal Oxide Semiconductor Field-effect Transistors. *Japanese journal of applied physics*, 43(4S):2151, 2004.
- [130] Ning Lu and Richard A Wachnik. Modeling of Resistance in FinFET Local Interconnect. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 62(8):1899–1907, 2015.
- [131] Ning Lu, Terence B Hook, Jeffrey B Johnson, Carl Wermer, Christopher Putnam, and Richard A Wachnik. Efficient and Accurate Schematic Transistor Model of FinFET Parasitic Elements. *IEEE Electron Device Letters*, 34(9):1100–1102, 2013.
- [132] Kai Ma and Xiaorui Wang. Pgcapping: exploiting power gating for power capping and core lifetime balancing in cmps. In *Proceedings of the 21st international conference on Parallel architectures and compilation techniques*, pages 13–22. ACM, 2012.
- [133] Paolo Madoglio, Hongtao Xu, Kailash Chandrashekar, Luis Cuellar, Muhammad Faisal, William Yee Li, Hyung Seok Kim, Khoa Minh Nguyen, Yulin Tan, Brent Carlton, et al. A 2.4 GHz WLAN digital polar transmitter with synthesized digital-totime converter in 14nm trigate/FinFET technology for IoT and wearable applications. In *Solid-State Circuits Conference (ISSCC), 2017 IEEE International*, pages 226–227. IEEE, 2017.
- [134] S Mahapatra, V Huard, A Kerber, V Reddy, S Kalpat, and A Haggag. Universality of nbti-from devices to circuits and products. In *Reliability Physics Symposium*, 2014 IEEE International, pages 3B–1. IEEE, 2014.

- [135] Souvik Mahapatra. Fundamentals of Bias Temperature Instability in MOS Transistors. Springer, 2016.
- [136] Souvik Mahapatra and Narendra Parihar. A review of nbti mechanisms and models. *Microelectronics Reliability*, 81:127–135, 2018.
- [137] W. P. Maszara and M. R. Lin. FinFETs Technology and Circuit Design Challenges. In 2013 Proceedings of the ESSCIRC (ESSCIRC), pages 3–8, Sept 2013. doi: 10. 1109/ESSCIRC.2013.6649058.
- [138] Witek P Maszara. Finfets: Designing for new logic technology. In Micro-and Nanoelectronics: Emerging Device Challenges and Solutions, pages 113–136. CRC Press, 2014.
- [139] Atmel SAM7SE microcontroller Education Kit:. http://www.atmel.com/tools/ sam7se-ek.aspx.
- [140] Evelyn Mintarno, Vishal Chandra, David Pietromonaco, Robert Aitken, and Robert W Dutton. Workload dependent nbti and pbti analysis for a sub-45nm commercial microprocessor. In *Reliability Physics Symposium (IRPS), 2013 IEEE International*, pages 3A–1. IEEE, 2013.
- [141] Evelyn Mintarno et al. Self-tuning for maximized lifetime energy-efficiency in the presence of circuit aging. *IEEE TCAD*, 30(5):760–773, 2011.
- [142] Prateek Mishra, Anish Muttreja, and Niraj K Jha. FinFET Circuit Design. In *Nano-electronic Circuit Design*, pages 23–54. Springer, 2011.
- [143] Subrat Mishra, Hiu Yung Wong, Ravi Tiwari, Ankush Chaudhary, Narendra Parihar, Rakesh Rao, Steve Motzny, Victor Moroz, and Souvik Mahapatra. Predictive tcad for nbti stress-recovery in various device architectures and channel materials. In *Reliability Physics Symposium (IRPS), 2017 IEEE International*, pages 6A–3. IEEE, 2017.
- [144] Subhasish Mitra, Kevin Brelsford, and Pia N Sanda. Cross-layer resilience challenges: Metrics and optimization. In *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2010, pages 1029–1034. IEEE, 2010.
- [145] Gordon E Moore. Cramming more components onto integrated circuits. *Proceedings* of the IEEE, 86(1):82–85, 1998.
- [146] Hassan Mostafa, Mohab Anis, and Mohamed Elmasry. Nbti and process variations compensation circuits using adaptive body bias. *Semiconductor Manufacturing, IEEE Transactions on*, 25(3):460–467, 2012.
- [147] Ann Mutschler. Transistor Aging Intensifies At 10/7nm And Below. https: //semiengineering.com/transistor-aging-intensifies-10nm/, 2017. [Online; accessed 13-July-2017].

- [148] Anish Muttreja, Niket Agarwal, and Niraj K Jha. CMOS Logic Design with Independent-gate FinFETs. In Computer Design, 2007. ICCD 2007. 25th International Conference on, pages 560–567. IEEE, 2007.
- [149] M Naouss and F Marc. Design and implementation of a low cost test bench to assess the reliability of fpga. *Microelectronics Reliability*, 55(9):1341–1345, 2015.
- [150] Katayoun Neshatpour, Wayne Burleson, Amin Khajeh, and Houman Homayoun. Enhancing Power, Performance, and Energy Efficiency in Chip Multiprocessors Exploiting Inverse Thermal Dependence. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 2018.
- [151] Hang Nguyen. Resiliency challenges in future communications infrastructure. In *IEEE Communications Quality and Reliability Workshop*, 2014.
- [152] Fabian Oboril and Mehdi B Tahoori. Extratime: A framework for exploration of clock and power gating for bti and hei aging mitigation. *ITG-Fachbericht-Zuverlässigkeit und Entwurf*, 2011.
- [153] Fabian Oboril and Mehdi B Tahoori. Reducing wearout in embedded processors using proactive fine-grain dynamic runtime adaptation. In *Test Symposium (ETS), 2012 17th IEEE European*, pages 1–6. IEEE, 2012.
- [154] Fabian Oboril and Mehdi B Tahoori. Cross-layer approaches for an aging-aware design of nanoscale microprocessors: Dissertation summary: Ieee tttc ej mccluskey doctoral thesis award competition finalist. In *Test Conference (ITC), 2015 IEEE International*, pages 1–10. IEEE, 2015.
- [155] Shin-ichi Ohfuji and Mitsuo Tsukada. Recovery of electric resistance degraded by electromigration. *Journal of applied physics*, 78(6):3769–3775, 1995.
- [156] Computing Community Consortium (CCC) Visioning Study on Cross-Layer Reliability:. http://www.relxlayer.org/.
- [157] S Pae, M Agostinelli, M Brazier, R Chau, G Dewey, T Ghani, M Hattendorf, J Hicks, J Kavalieros, K Kuhn, et al. Bti reliability of 45 nm high-k+ metal-gate process technology. In *Reliability Physics Symposium, 2008. IRPS 2008. IEEE International*, pages 352–357. IEEE, 2008.
- [158] Narendra Parihar, Uma Sharma, Subhadeep Mukhopadhyay, Nilesh Goel, Ankush Chaudhary, Rakesh Rao, and Souvik Mahapatra. Resolution of disputes concerning the physical mechanism and DC/AC stress/recovery modeling of Negative Bias Temperature Instability (NBTI) in p-MOSFETs. In *Reliability Physics Symposium (IRPS)*, 2017 IEEE International, pages XT–1. IEEE, 2017.
- [159] Bipul C Paul, Kunhyuk Kang, Haldun Kufluoglu, Muhammad Alam, Kaushik Roy, et al. Impact of nbti on the temporal performance degradation of digital circuits. *Electron Device Letters, IEEE*, 26(8):560–562, 2005.

- [160] A. Paya and D. Marinescu. Energy-aware load balancing and application scaling for the cloud ecosystem. *IEEE Transactions on Cloud Computing*, PP(99):1–1, 2015. ISSN 2168-7161. doi: 10.1109/TCC.2015.2396059.
- [161] Nathaniel Pinckney, Lucian Shifren, Brian Cline, Saurabh Sinha, Supreet Jeloka, Ronald G Dreslinski, Trevor Mudge, Dennis Sylvester, and David Blaauw. Nearthreshold Computing in FinFET Technologies: Opportunities for Improved Voltage Scalability. In *Proceedings of the 53rd Annual Design Automation Conference*, page 76. ACM, 2016.
- [162] Gregor Pobegen, Thomas Aichinger, Michael Nelhiebel, and Tibor Grasser. Understanding temperature acceleration for nbti. In *Proc. Intl. Electron Devices Meeting* (*IEDM*), pages 27–3, 2011.
- [163] Gracieli Posser, Sachin S Sapatnekar, and Ricardo Reis. Analyzing the electromigration effects on different metal layers and different wire lengths. In *Electromigration Inside Logic Cells*, pages 93–98. Springer, 2017.
- [164] Cadence Voltus IC Power Integrity Solution: Rapid power signoff and design closure. https://www.cadence.com/content/cadence-www/global/en_US/home/tools/ digital-design-and-signoff/silicon-signoff/voltus-ic-power-integrity-solution.html.
- [165] C Prasad, S Ramey, and L Jiang. Self-heating in advanced cmos technologies. In *Reliability Physics Symposium (IRPS), 2017 IEEE International*, pages 6A–4. IEEE, 2017.
- [166] Zhenyu Qi and Mircea R Stan. Nbti resilient circuits using adaptive body biasing. In Proceedings of the 18th ACM Great Lakes symposium on VLSI, pages 285–290. ACM, 2008.
- [167] Jan M Rabaey, Anantha P Chandrakasan, and Borivoje Nikolic. *Digital integrated circuits*, volume 2. Prentice hall Englewood Cliffs, 2002.
- [168] KK Ramakrishnan, Smitha Suresh, Narayanan Vijaykrishnan, Mary Jane Irwin, and Vijay Degalahal. Impact of nbti on fpgas. In VLSI Design, 2007. Held jointly with 6th International Conference on Embedded Systems., 20th International Conference on, pages 717–722. IEEE, 2007.
- [169] SM Ramey, C Prasad, and A Rahman. Technology scaling implications for BTI reliability. *Microelectronics Reliability*, 82:42–50, 2018.
- [170] M Rashed, N Jain, J Kim, M Tarabbia, I Rahim, S Ahmed, Je Kim, I Lin, S Chan, H Yoshida, et al. Innovations in Special Constructs for Standard Cell Libraries in Sub 28nm Technologies. In *Electron Devices Meeting (IEDM)*, 2013 IEEE International, pages 9–7. IEEE, 2013.
- [171] Ricardo Reis, Yu Cao, and Gilson Wirth. Circuit design for reliability. Springer, 2015.

- [172] Hans Reisinger, Oliver Blank, Wolfgang Heinrigs, Wolfgang Gustin, and Christian Schlünder. A comparison of very fast to very slow components in degradation and recovery due to nbti and bulk hole trapping to existing physical models. *Device and Materials Reliability, IEEE Transactions on*, 7(1):119–129, 2007.
- [173] ITRS Report: http://www.itrs2.net/itrs-reports.html.
- [174] IJ Ringler and JR Lloyd. Stress relaxation in pulsed dc electromigration measurements. *AIP Advances*, 6(9):095118, 2016.
- [175] Silvestre Salas Rodriguez, Julio C Tinoco, Andrea G Martinez-Lopez, Joaquín Alvarado, and Jean-Pierre Raskin. Parasitic Gate Capacitance Model for Triple-gate FinFETs. *IEEE Transactions on Electron Devices*, 60(11):3710–3717, 2013.
- [176] Daniele Rossi, Vasileios Tenentes, Bashir Al-Hashimi, et al. Nbti and leakage aware sleep transistor design for reliable and energy efficient power gating. *Proceedings of the IEEE European Test Symposium*, 2015.
- [177] Kaushik Roy, Saibal Mukhopadhyay, and Hamid Mahmoodi-Meimand. Leakage current mechanisms and leakage reduction techniques in deep-submicrometer cmos circuits. *Proceedings of the IEEE*, 91(2):305–327, 2003.
- [178] Pablo Royer, Paul Zuber, Binjie Cheng, Asen Asenov, and Marisa Lopez-Vallejo. Circuit-level Modeling of FinFET Sub-threshold Slope and DIBL Mismatch Beyond 22nm. In Simulation of Semiconductor Processes and Devices (SISPAD), 2013 International Conference on, pages 204–207. IEEE, 2013.
- [179] Takayasu Sakurai and A Richard Newton. Alpha-Power Law MOSFET Model and its Applications to CMOS Inverter Delay and Other Formulas. *IEEE Journal of Solid-State Circuits*, 25(2):584–594, 1990.
- [180] S Sarma, N Dutt, N Venkatasubramanian, A Nicolau, and P Gupta. Cyberphysical system-on-chip (cpsoc): Sensoractuator rich self-aware computational platform. University of California Irvine, Tech. Rep. CECS TR-13-06, 2013.
- [181] Jeff Sather. Battery technologies for iot. In *Enabling the Internet of Things*, pages 409–440. Springer, 2017.
- [182] Yasuo Sato, Masafumi Monden, Yousuke Miyake, and Seiji Kajihara. Reduction of nbti-induced degradation on ring oscillators in fpga. In *Dependable Computing* (*PRDC*), 2014 IEEE 20th Pacific Rim International Symposium on, pages 59–67. IEEE, 2014.
- [183] Christian Schlünder et al. On the influence of BTI and HCI on parameter variability. In *Reliability Physics Symposium (IRPS), 2017 IEEE International*, pages 2E–4. IEEE, 2017.

- [184] D. C. Sekar et al. Electromigration Resistant Power Delivery Systems. *IEEE Electron Device Letters*, 28(8):767–769, Aug 2007. ISSN 0741-3106. doi: 10.1109/LED.2007. 902165.
- [185] Deepak C Sekar, Bing Dang, Jeffrey A Davis, and James D Meindl. Electromigration resistant power delivery systems. *IEEE electron device letters*, 28(8):767–769, 2007.
- [186] Dipak Sengupta and Sachin S Sapatnekar. Predicting circuit aging using ring oscillators. In *Design Automation Conference (ASP-DAC), 2014 19th Asia and South Pacific,* pages 430–435. IEEE, 2014.
- [187] Dimitrios Serpanos and Marilyn Wolf. IoT Devices. In *Internet-of-Things (IoT)* Systems, pages 17–23. Springer, 2018.
- [188] Muhammad Shafique, Siddharth Garg, Jörg Henkel, and Diana Marculescu. The eda challenges in the dark silicon era: Temperature, reliability, and variability perspectives. In *Proceedings of the 51st Annual Design Automation Conference*, pages 1–6. ACM, 2014.
- [189] Nimay Shah, Rupak Samanta, Ming Zhang, Jiang Hu, and Duncan Walker. Built-in proactive tuning system for circuit aging resilience. In *Defect and Fault Tolerance of VLSI Systems, 2008. DFTVS'08. IEEE International Symposium on*, pages 96–104. IEEE, 2008.
- [190] Farhana Sheikh and Vidya Varadarajan. The Impact of Device-width Quantization on Digital Circuit Design Using FinFET Structures. *Proc. EE241 Spring*, 1, 2004.
- [191] Jeonghee Shin, Victor Zyuban, Pradip Bose, and Timothy M Pinkston. A proactive wearout recovery approach for exploiting microarchitectural redundancy to extend cache sram lifetime. In ACM SIGARCH Computer Architecture News, volume 36, pages 353–362. IEEE Computer Society, 2008.
- [192] Taniya Siddiqua and Sudhanva Gurumurthi. Nbti-aware dynamic instruction scheduling. In *Proceedings of the 5th Workshop on Silicon Errors in Logic-System Effects*. Citeseer, 2009.
- [193] Taniya Siddiqua and Sudhanva Gurumurthi. Recovery boosting: A technique to enhance nbti recovery in sram arrays. In *VLSI (ISVLSI), 2010 IEEE Computer Society Annual Symposium on*, pages 393–398. IEEE, 2010.
- [194] Prashant Singh, Eric Karl, David Blaauw, and Dennis Sylvester. Compact degradation sensors for monitoring nbti and oxide degradation. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 20(9):1645–1655, 2012.
- [195] Ajith Sivadasan, Florian Cacho, Sidi Ahmed Benhassain, Vincent Huard, and Lorena Anghel. Study of workload impact on bti hci induced aging of digital circuits. In *Proceedings of the 2016 Conference on Design, Automation & Test in Europe*, pages 1020–1021. EDA Consortium, 2016.

- [196] Ajith Sivadasan, S Mhira, Armelle Notin, A Benhassain, V Huard, Etienne Maurin, F Cacho, L Anghel, and A Bravaix. Architecture-and workload-dependent digital failure rate. In *Reliability Physics Symposium (IRPS), 2017 IEEE International*, pages CR-8. IEEE, 2017.
- [197] T. Song, W. Rim, S. Park, Y. Kim, G. Yang, H. Kim, S. Baek, J. Jung, B. Kwon, S. Cho, H. Jung, Y. Choo, and J. Choi. A 10 nm finfet 128 mb sram with assist adjustment system for power, performance, and area optimization. *IEEE Journal of Solid-State Circuits*, 52(1):240–249, Jan 2017. ISSN 0018-9200. doi: 10.1109/JSSC. 2016.2609386.
- [198] Warin Sootkaneung, Sasithorn Chookaew, and Suppachai Howimanporn. Combined Impact of BTI and Temperature Effect Inversion on Circuit Performance. In Asian Test Symposium (ATS), 2016 IEEE 25th, pages 310–315. IEEE, 2016.
- [199] Jayanth Srinivasan, Sarita V Adve, Pradip Bose, and Jude A Rivers. The case for lifetime reliability-aware microprocessors. In ACM SIGARCH Computer Architecture News, volume 32, page 276. IEEE Computer Society, 2004.
- [200] Jayanth Srinivasan, Sarita V Adve, Pradip Bose, and Jude A Rivers. Exploiting structural duplication for lifetime reliability enhancement. In *Computer Architecture*, 2005. ISCA'05. Proceedings. 32nd International Symposium on, pages 520–531. IEEE, 2005.
- [201] Mircea R Stan and Paolo Re. Electromigration-aware design. In *Circuit Theory and Design, 2009. ECCTD 2009. European Conference on*, pages 786–789. IEEE, 2009.
- [202] James H Stathis, M Wang, RG Southwick, EY Wu, BP Linder, EG Liniger, G Bonilla, and H Kothari. Reliability challenges for the 10nm node and beyond. In *Electron Devices Meeting (IEDM), 2014 IEEE International*, pages 20–6. IEEE, 2014.
- [203] James H Stathis, Souvik Mahapatra, and Tibor Grasser. Controversial issues in negative bias temperature instability. *Microelectronics Reliability*, 81:244–251, 2018.
- [204] V Sukharev, X Huang, and SX-D Tan. Electromigration induced stress evolution under alternate current and pulse current loads. *Journal of Applied Physics*, 118(3): 034504, 2015.
- [205] Vikram B Suresh and Wayne P Burleson. Fine grained wearout sensing using metastability resolution time. In *Quality Electronic Design (ISQED), 2014 15th International Symposium on*, pages 480–485. IEEE, 2014.
- [206] Brian Swahn and Soha Hassoun. Gate Sizing: FinFETs vs 32nm Bulk MOSFETs. In Design Automation Conference, 2006 43rd ACM/IEEE, pages 528–531. IEEE, 2006.
- [207] Shiva Taghipour and Rahebeh Niaraki Asli. Aging Comparative Analysis of Highperformance FinFET and CMOS Flip-flops. *Microelectronics Reliability*, 69:52–59, 2017.

- [208] Sheldon X-D Tan, Hussam Amrouch, Taeyoung Kim, Zeyu Sun, Chase Cook, and Jörg Henkel. Recent advances in em and bti induced reliability modeling, analysis and optimization. *Integration, the VLSI Journal*, 2017.
- [209] Jiang Tao, Jone F Chen, Nathan W Cheung, and Chenming Hu. Modeling and characterization of electromigration failures under bidirectional current stress. *Electron Devices, IEEE Transactions on*, 43(5):800–808, 1996.
- [210] J Teshima and Jamil J Clarke. From transistors to bumps: Preparing sem crosssections by combining site-specific cleaving and broad ion beam milling. SOLID STATE TECHNOLOGY, 58(7):21–26, 2014.
- [211] Mastering the Magic of Multi-Patterning:. http://go.mentor.com/4gue4.
- [212] Abhishek Tiwari and Josep Torrellas. Facelift: Hiding and slowing down aging in multicores. In *Microarchitecture*, 2008. MICRO-41. 2008 41st IEEE/ACM International Symposium on, pages 129–140. IEEE, 2008.
- [213] James W Tschanz, James T Kao, Siva G Narendra, Raj Nair, Dimitri A Antoniadis, Anantha P Chandrakasan, and Vivek De. Adaptive Body Bias for Reducing Impacts of Die-to-die and Within-die Parameter Variations on Microprocessor Frequency and Leakage. *IEEE Journal of Solid-State Circuits*, 37(11):1396–1402, 2002.
- [214] Spencer Tu. Putting the Pieces Together in the Materials Space: Advanced Materials Solutions for 10nm and Beyond. In *SEMICON Taiwan*, 2015.
- [215] Victor van Santen et al. Reliability in Super-and Near-Threshold Computing: A Unified Model of RTN, BTI, and PV. *TCAS-I*, 2017.
- [216] Rakesh Vattikonda, Wenping Wang, and Yu Cao. Modeling and minimization of pmos nbti effect for robust nanometer design. In *Proceedings of the 43rd annual Design Automation Conference*, pages 1047–1052. ACM, 2006.
- [217] Jyothi Bhaskarr Velamala, Ketul Sutaria, Takashi Sato, and Yu Cao. Physics matters: statistical aging prediction under trapping/detrapping. In *Proceedings of the 49th Annual Design Automation Conference*, pages 139–144. ACM, 2012.
- [218] Jyothi Bhaskarr Velamala, Ketul B Sutaria, Hirofumi Shimizu, Hiromitsu Awano, Takashi Sato, Gilson Wirth, and Yu Cao. Compact modeling of statistical bti under trapping/detrapping. *IEEE Transactions on Electron Devices*, 60(11):3645–3654, 2013.
- [219] Siva Velusamy, Wei Huang, John Lach, Mircea Stan, and Kevin Skadron. Monitoring temperature in fpga based socs. In *Computer Design: VLSI in Computers and Processors, 2005. ICCD 2005. Proceedings. 2005 IEEE International Conference on*, pages 634–637. IEEE, 2005.

- [220] Alice Wang, Anantha P Chandrakasan, and Stephen V Kosonocky. Optimal Supply and Threshold Scaling for Subthreshold CMOS Circuits. In *VLSI*, 2002. Proceedings. *IEEE Computer Society Annual Symposium on*, pages 7–11. IEEE, 2002.
- [221] Runsheng Wang, Pengpeng Ren, Changze Liu, Shaofeng Guo, and Ru Huang. Understanding nbti-induced dynamic variability in the nano-reliability era: From devices to circuits. In *Physical and Failure Analysis of Integrated Circuits (IPFA), 2015 IEEE* 22nd International Symposium on the, pages 119–121. IEEE, 2015.
- [222] S. Wang, T. Kim, Z. Sun, S. X. D. Tan, and M. B. Tahoori. Recovery-Aware Proactive TSV Repair for Electromigration Lifetime Enhancement in 3-D ICs. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, PP(99):1–13, 2017. ISSN 1063-8210. doi: 10.1109/TVLSI.2017.2775586.
- [223] Wenping Wang, Shengqi Yang, Sarvesh Bhardwaj, Rakesh Vattikonda, Sarma Vrudhula, Frank Liu, and Yu Cao. The impact of nbti on the performance of combinational and sequential circuits. In *Proceedings of the 44th annual Design Automation Conference*, pages 364–369. ACM, 2007.
- [224] Xiaoyang Wang, Po-Han Peter Wang, Yuan Cao, and Patrick P Mercier. A 0.6 v 75nw all-cmos temperature sensor with 1.67 m° c/mv supply sensitivity. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 2017.
- [225] Xingsheng Wang, Andrew R Brown, Binjie Cheng, and Asen Asenov. Statistical Variability and Reliability in Nanoscale FinFETs. In *Electron Devices Meeting (IEDM)*, 2011 IEEE International, pages 5–4. IEEE, 2011.
- [226] James Warnock. Circuit Design Challenges at the 14nm Technology Node. In *Proceedings of the 48th Design Automation Conference*, pages 464–467. ACM, 2011.
- [227] O Weber. FDSOI vs FinFET: differentiating device features for ultra low power & IoT applications. In *IC Design and Technology (ICICDT), 2017 IEEE International Conference on*, pages 1–3. IEEE, 2017.
- [228] Piotr Weber, Maciej Zagrabski, Przemyslaw Musz, Krzysztof Kepa, Maciej Nikodem, and Bartosz Wojciechowski. Configurable heat generators for fpgas. In *Thermal Investigations of ICs and Systems (THERMINIC), 2014 20th International Workshop on*, pages 1–4. IEEE, 2014.
- [229] Neil HE Weste and David Money Harris. *CMOS VLSI Design: A Circuits and Systems Perspective*. Pearson Addison-Wesley, 2005.
- [230] David Wolpert and Paul Ampadu. Temperature effects in semiconductors. In *Managing* temperature effects in nanoscale adaptive systems, pages 15–33. Springer, 2012.
- [231] Stuart N Wooters, Adam C Cabe, Zhenyu Qi, Jiajing Wang, Randy W Mann, Benton H Calhoun, Mircea R Stan, and Travis N Blalock. Tracking on-chip age using distributed,
embedded sensors. Very Large Scale Integration (VLSI) Systems, IEEE Transactions on, 20(11):1974–1985, 2012.

- [232] Kai-Chiang Wu, Diana Marculescu, Ming-Chao Lee, and Shih-Chieh Chang. Analysis and mitigation of nbti-induced performance degradation for power-gated circuits. In *Proceedings of the 17th IEEE/ACM international symposium on Low-power electronics and design*, pages 139–144. IEEE Press, 2011.
- [233] Kai-Chiang Wu, Chao Lin, Yao-Te Wang, and Shuen-Shiang Yang. Bti-aware sleep transistor sizing algorithm for reliable power gating designs. *Computer-Aided Design* of Integrated Circuits and Systems, IEEE Transactions on, 33(10):1591–1595, 2014.
- [234] Kaiyuan Yang, Qing Dong, Wanyeong Jung, Yiqun Zhang, Myungjoon Choi, David Blaauw, and Dennis Sylvester. 9.2 a 0.6 nj- 0.22/+ 0.19° c inaccuracy temperature sensor using exponential subthreshold oscillation dependence. In *Solid-State Circuits Conference (ISSCC), 2017 IEEE International*, pages 160–161. IEEE, 2017.
- [235] Wen-Kuan Yeh, Wenqi Zhang, Yi-Lin Yang, An-Ni Dai, Kehuey Wu, Tung-Huan Chou, Cheng-Li Lin, Kwang-Jow Gan, Chia-Hung Shih, and Po-Ying Chen. The Observation of Width Quantization Impact on Device Performance and Reliability for High-k/Metal Tri-Gate FinFET. *IEEE Transactions on Device and Materials Reliability*, 16(4):610–616, 2016.
- [236] Bin Yu, Leland Chang, Shibly Ahmed, Haihong Wang, Scott Bell, Chih-Yuh Yang, Cyrus Tabery, Chau Ho, Qi Xiang, Tsu-Jae King, et al. FinFET Scaling to 10 nm Gate Length. In *Electron Devices Meeting*, 2002. *IEDM'02. International*, pages 251–254. IEEE, 2002.
- [237] S Zafar, YH Kim, V Narayanan, C Cabral Jr, V Paruchuri, B Doris, J Stathis, A Callegari, and M Chudzik. A comparative study of nbti and pbti (charge trapping) in sio2/hfo2 stacks with fusi, tin, re gates. In VLSI Technology, 2006. Digest of Technical Papers. 2006 Symposium on, pages 23–25. IEEE, 2006.
- [238] Lide Zhang and Robert P Dick. Scheduled voltage scaling for increasing lifetime in the presence of nbti. In *Proceedings of the Asia and South Pacific Design Automation Conference*, pages 492–497. IEEE Press, 2009.
- [239] Xuehui Zhang, Kan Xiao, and Mohammad Tehranipoor. Path-delay fingerprinting for identification of recovered ics. In *Defect and Fault Tolerance in VLSI and Nanotechnology Systems (DFT), 2012 IEEE International Symposium on*, pages 13–18. IEEE, 2012.
- [240] Huajie Zhou, Yi Song, Qiuxia Xu, Yongliang Li, and Huaxiang Yin. Fabrication of Bulk-Si FinFET Using CMOS Compatible Process. *Microelectronic Engineering*, 94: 26–28, 2012.

- [241] Cheng Zhuo, Kaviraj Chopra, Dennis Sylvester, and David Blaauw. Process variation and temperature-aware full chip oxide breakdown reliability analysis. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 30(9):1321–1334, 2011.
- [242] B. Zimmer, S. O. Toh, H. Vo, Y. Lee, O. Thomas, K. Asanovic, and B. Nikolic. Sram assist techniques for operation in a wide voltage range in 28-nm cmos. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 59(12):853–857, Dec 2012. ISSN 1549-7747. doi: 10.1109/TCSII.2012.2231015.
- [243] Yazhou Zu, Wei Huang, Indrani Paul, and Vijay Janapa Reddi. T_i-states: Processor Power Management in the Temperature Inversion Region. In *Microarchitecture* (*MICRO*), 2016 49th Annual IEEE/ACM International Symposium on, pages 1–13. IEEE, 2016.