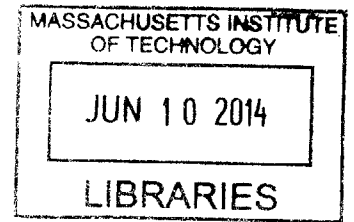


**Energy-Aware System Design Using Circuit
Reconfigurability with a Focus on Low-Power
SRAMs**

ARCHIVES

by

Yildiz Sinangil



B.S. in Electrical Engineering, Bogazici University, 2008
S.M. in Electrical Engineering, Massachusetts Institute of Technology,
2010

Submitted to the Department of Electrical Engineering and Computer
Science

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2014

© Massachusetts Institute of Technology 2014. All rights reserved.

Signature redacted

Author
Department of Electrical Engineering and Computer Science

May 21, 2014

Signature redacted

Certified by
Anantha P. Chandrakasan

Joseph F. and Nancy P. Keithley Professor of Electrical Engineering
Thesis Supervisor

Signature redacted

Accepted by
Leslie A. Kolodziejcki

Chairman, Department Committee on Graduate Students

Energy-Aware System Design Using Circuit Reconfigurability with a Focus on Low-Power SRAMs

by

Yildiz Sinangil

Submitted to the Department of Electrical Engineering and Computer Science
on May 21, 2014, in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy

Abstract

Today's complex systems generally target competing design goals such as maximizing performance while minimizing energy. Moreover, they have to work efficiently under changing system dynamics and application loads. Thus, for better power and performance optimization, they need to adapt to different conditions on-the-fly. In this regard, systems need to monitor important metrics such as energy consumption and performance.

First part of this thesis focuses on an energy monitoring circuit design that can generate a digital representation of the absolute energy per operation of a circuit. A test-chip is fabricated in a 65nm LP CMOS process and the energy monitoring circuit is demonstrated for an SRAM application. The small power ($< 0.1\%$) and area overhead (16%) of the energy monitor motivated its usage within a system level application.

Next, as a collaboration of circuit designers and architects, energy monitoring is extended to a processor system for system-level power and performance optimizations. To enable self-adaptation, custom blocks in the system are designed to accommodate reconfigurability. In this regard, the data cache is designed to be reconfigurable in terms of set associativity (1-4 sets) and size (1kB to 4kB per set). Furthermore, the system is designed to enable voltage and frequency scaling. Energy monitoring circuits that are capable of monitoring runtime conditions for different domains are embedded into on-chip DC-DC converters. Those ideas are demonstrated in a $0.18\mu\text{m}$ system-on-chip design and measurement results show that the system can achieve up to $8.4\times$ energy savings.

SRAMs account for a large fraction of system area and power. Thus, low-voltage SRAM operation is crucial for an energy-efficient system design. Two SRAM circuits are demonstrated in 65nm and $0.18\mu\text{m}$ test-chips using 8T bit-cells and write assists, and they are measured to be voltage scalable from 1V to 0.37V and from 1.8V to 0.6V respectively. Secondly, write and read assist techniques are illustrated for an industry sized, 6T bit-cell based SRAM design. A test-chip is fabricated using a cutting-edge 28nm FD-SOI technology and the SRAMs are measured to be operational down to 0.43 V.

Thesis Supervisor: Anantha P. Chandrakasan

Title: Joseph F. and Nancy P. Keithley Professor of Electrical Engineering

Acknowledgments

First and foremost, I would like to thank Professor Anantha Chandrakasan for his guidance and mentorship. Throughout these years, I learned a lot from you; not only on technical issues, but also on topics of how to be a good, fair and respectful researcher. Thank you for being a strong supporter of my work, and always making yourself available for discussion despite your busy schedule. I would also like to extend my sincerest gratitudes to my PhD. thesis committee members for their contributions and support: Prof. David Perreault and Prof. Saman Amarasinghe. It was an honor to share my work with you. Thank you for widening the perspective of this thesis with your invaluable feedback.

There are several faculty members at MIT who have had a profound impact on my life at MIT, both technically and non-technically. I would like to thank Prof. Srin Devadas and Prof. Anant Agarwal for their support and collaboration throughout the Angstrom and PERFECT projects. I am grateful to Prof. Jeffrey Lang, Prof. David Perreault, Prof. John Kassakian and Prof. Khurram Afridi for sharing their expertise in teaching with me. I would like to thank Prof. Dina Katabi for our encouraging discussions about how to be a successful woman in engineering.

Being a member of the ananthagroup was a privilege since I had the opportunity to meet and work with a large number of smart and fun people. I would like to start with thanking Daniel Finchelstein for being a great mentor during my summer internship. I remember the summer time when he was teaching me the basics about Verilog and FPGA programming in front of the white-board. Thank you for pushing me to my limits, making me learn tremendously, and helping me to make the summer a great success. I also had the opportunity to have lunch-time discussions with Naveen, Yogesh, Joyce, Vivienne, Payam, and Denis. It was great to get to know you even before coming to MIT as a graduate student. Some other past graduate students I was fortunate enough to overlap were Marcus, Masood, Fred, Nigel, Jose, Brian, Tao, Hye Won, and Courtney. I am grateful to get to know you.

I am privileged to work and collaborate with current students from ananthagroup

and other groups. First I would like to acknowledge the SEEC test-chip group: Sabrina, Mahmut, Nathan, George, Eric, Jason, Hank, and all the UROP students that contributed to the project. We spent many hours in our meetings and the self-aware test-chip would not be possible without your contributions. I would like to thank Avishek and Chu for our discussions about SRAM circuits, switched capacitors and prediction algorithms. It has been a great experience to work with you and I learned a lot throughout the process. I would like to especially thank Nathan for not only our collaborations during the tape-outs but also for his tremendous help in digital flow, board design and testing.

The most rewarding experience of MIT was definitely being surrounded by lots of smart, talented and fun people. Bonnie - Thanks for being my ananthagroup twin for the last six years. I had great fun having you sitting at the next cubicle throughout my two internships at Texas Instruments. Rahul and Saurav - Thank you for sharing the last six years with me. It was great to talk to you about not only technical work, but also about all the deadlines, job search hassles and tape-out tips. Arun - Thanks for our discussions about TSMC tape-out specifics. Frank and Michael - Thank you very much for your help in moving our apartment from Medford to Burlington! I would also like to thank Dina, Sungjae, Omid, Georgios, Nachiket, Rui, Chiraag, Priyanka, Mehul, Gilad, Theresa and many others for being wonderful friends and colleagues. I would like to thank Phillip, Farnaz and Katherine for sharing the MTL Annual Research Conference with me. It was a pleasure to be in the committee. Lastly, I would like to thank Margaret for helping the lab run smoothly. Seeing her smiling face every day was a blessing. I would like to thank Patrick and Chelce for being great neighbors and the wonderful movie and game nights. Thanks to Shweta for our fun talks. Lastly, I would like to thank the Turkish gang - Gozde, Filiz, Renin, Oguzhan, Cagatay and others for our great times together.

This work would not be possible without the endless support and love of my family. To my parents, Nesrin and Necati, Havva and Sukru and my brother Gunes and sister Tugba: Thank you for being there for me. Although you are miles away, you were always supporting me with your frequent phone calls. And finally I would

like to thank my dear husband: Mahmut Ersin Sinangil. This thesis would not be possible without your endless support and encouragement. You have been the one to set me straight when I was stressed out and frustrated and you were the one to share my accomplishments and successes. Thanks for being my best friend, my soul mate and my best companion.

Contents

1	Introduction	27
1.1	Self-Aware Systems	30
1.2	Energy-efficiency and voltage-scaling	33
1.3	Low-Voltage SRAM Design	35
1.3.1	Sources of Variation	36
1.3.2	SRAM Failure Mechanisms	37
1.4	Thesis Contributions	44
2	Energy Sensing Circuit for an SRAM Application	49
2.1	Overview of Energy Sensing Concept	51
2.2	Energy Monitoring Circuit Challenges	54
2.2.1	Effect of ΔV Drop on SRAM Operation	54
2.2.2	Effect of ΔV Selection On Accuracy	56
2.2.3	Effect of Comparator Offset on Accuracy	57
2.2.4	Effect of Non-Idealities of the Capacitor	58
2.3	Energy Sensor Demonstration for an SRAM Application	60
2.4	Measurement Results	68
2.4.1	Measured EOP Accuracy	70
2.4.2	EOP Under Dynamic Effects	71
2.5	Summary and Conclusions	72
3	Self-Aware Processor Design Using Embedded Energy Monitors	75
3.1	Observe-Decide-Act Loop For Self-Aware Processor	77

3.2	LEON3 Processor	79
3.3	The Self-Aware Processor System-on-Chip Implementation	80
3.4	DC-DC converters with Embedded Energy Monitors	83
3.4.1	DC-DC Converter Implementation for the Self-Aware Processor	83
3.4.2	Energy Monitoring Circuit Operation with DC-DC	88
3.4.3	Results for the DC-DC Converters with Embedded Energy Monitors	92
3.5	Reconfigurable SRAM Design	95
3.5.1	Effect of Set Associativity and Size of L1-Cache	97
3.5.2	Single Cycle Tag Invalidation	99
3.6	Measurement Results of the Self-Aware Processor	100
3.6.1	Effect of DVFS and D-Cache Adaptation to Energy Consumption	104
3.6.2	Effect of D-Cache Adaptation for Different Applications	105
3.6.3	Comparison to a Conventional Race-to-Idle System	107
3.6.4	Comparison of the Self-Aware System to Recent Work	108
3.7	Summary and Conclusions	109
4	Low-Voltage SRAM Design	113
4.1	8T Bit-cell Based SRAM Design	114
4.1.1	8T and 6T bit-cell Topologies	115
4.1.2	Error Map Comparison for 6T vs. 8T	117
4.1.3	Write-Assist Circuit	118
4.1.4	Memory Organization of the 128kb SRAM	121
4.1.5	Two-stage Sensing	123
4.1.6	Sense-Amplifier Offset Reduction Using Body-Biasing	125
4.1.7	Measurement Results	130
4.2	A 6T Bit-cell Based SRAM in 28nm FD-SOI CMOS for Operation Below 0.5 V	131
4.2.1	FD-SOI Technology Specifics	132
4.2.2	Low-Voltage Challenges of 6T Bit-cell Operation in FD-SOI	135

4.2.3	Write-Assist using Body-Biasing	137
4.2.4	Memory Building Blocks	143
4.2.5	Improving Read Stability Using Hierarchical Bit-lines	150
4.2.6	Test-chip Results	151
4.3	Summary and Conclusions	156
5	Conclusions and Future Work	161
5.1	Summary of Contributions	162
5.1.1	Energy Monitoring Circuit Demonstration for an SRAM Application	162
5.1.2	Self-Aware System That can Reconfigure Itself Based On Actual Energy Measurements	163
5.1.3	Low-Voltage SRAM Design	164
5.2	Future Work	165
A	Yield Calculation of SRAMs	169
B	SRAM Assist Techniques	171
C	LEON3 Processor	173
D	Switching Regulators	177
D.0.1	Conduction Loss	180
D.0.2	Switching Loss	180
D.0.3	Timing Losses	181
D.0.4	Leakage Loss	181
D.0.5	Efficiency of the PFM Mode DC-DC Converter	181

List of Figures

1-1	Li-ion battery energy density over the years.	28
1-2	Chip complexity over the last two decades.	28
1-3	Trend of cache integration over the last two decades.	29
1-4	Energy and frequency scaling with voltage for a microprocessor application.	34
1-5	(a) Itanium microprocessor with 54MB [1], and (b) Power8 microprocessor with 96MB [2] on-chip caches.	35
1-6	Bit-cell area scaling over technology nodes. [3]	36
1-7	SRAM devices V_T variation scaling trend [4].	36
1-8	SRAM failures due to hard errors and stability issues versus technology nodes [4].	38
1-9	6T bit-cell in read operation.	39
1-10	RSNM butterfly curves for a 28nm bit-cell for 100 MC runs.	40
1-11	A typical SRAM read path.	41
1-12	6T bit-cell in write operation.	42
1-13	Illustration of various peripheral assist techniques. A similar diagram is illustrated in [5].	44
2-1	SandyBridge power management architecture [6].	50
2-2	(a) Normal operation, (b) Energy sensing cycle, and (c) Voltage recovery.	53
2-3	Selection of ΔV at low voltages for cell stability.	55
2-4	The required frequency adjustment for a $\Delta V=10\text{mV}$	56
2-5	Selection of ΔV for accuracy considerations.	57

2-6	Using aluminum capacitor with a large ESR (8Ω) results into significant error (10%) in energy measurement whereas MLCC parasitics bring negligible error.	59
2-7	Chip block diagram with a focus on the architecture of the energy-sensing circuit.	60
2-8	Schematics of the DAC.	62
2-9	DAC output span for 128 different input words.	63
2-10	A 4 bit ripple-borrow subtractor using full subtractors.	64
2-11	Block diagram of the custom design shift-and-add multiplier.	65
2-12	Block diagram of the custom design shift-and-subtract divider.	66
2-13	Trimmed strong-arm type sense amplifier is used for comparator.	67
2-14	Die photo of the 128 kb SRAM in 65 nm CMOS.	67
2-15	Oscilloscope outputs for the critical signals of energy sensing circuit.	69
2-16	EOP vs. V_{DD} graph using three methods: 1- measured, 2- simulated, and 3- energy monitoring circuit output.	71
2-17	Measured EOP under different read operation to write operation ratios and temperatures.	72
3-1	Angstrom multicore architecture.	76
3-2	ODA loop for the self-aware processor.	77
3-3	LEON3 processor core block diagram.	79
3-4	The self-aware processor block diagram.	80
3-5	Hardware and software separation in self-aware processor system.	82
3-6	The variation of the load power vs. V_{DD}	84
3-7	The block diagram of the DC-DC converters with the embedded energy monitoring circuit.	85
3-8	PMOS pulse (T_{PMOS}) generation circuit.	86
3-9	Ratio of required PMOS pulse width to NMOS pulse width vs. operating voltage.	87
3-10	The operation of the energy monitors and DC-DC converters together.	88

3-11	The asynchronous boundary block diagram.	89
3-12	The asynchronous boundary signals between the fixed domain and the core domain.	90
3-13	The reference voltage levels for energy sensing and voltage change operations.	92
3-14	The superior efficiency curve can be selected for better efficiency as shown in the simulation result.	93
3-15	The oscilloscope output of the system while performing energy sensing and voltage change operations.	95
3-16	Adaptive d-cache structure.	96
3-17	Graphite simulation for changing cache size and associativity on a single core system. Different applications require different optimizations for minimum energy.	98
3-18	Single-cycle tag invalidation.	99
3-19	Die photo of the 0.18 μ m self-aware processor system.	101
3-20	Area breakdown of the test-chip prototype.	102
3-21	Block diagram of the test setup of the self-aware processor system. . .	103
3-22	Picture of the test setup of the self-aware processor system.	104
3-23	(a) The effect of DVFS on energy consumption, and (b) the effect of DVFS and cache adaptation.	105
3-24	The effect of cache adaptation on energy consumption.	106
3-25	System simulation running four phases of a multi-media application using 1-self-aware adaptation, 2-static configuration with race-to-idle operation.	107
4-1	Bit-cell voltage scaling over technology nodes [3].	114
4-2	Schematics of (a) a 6T bit-cell, and (b) an 8T bit-cell.	115
4-3	Error map simulation of a 64 \times 128 SRAM block designed using a (a) 6T bit-cell, (b) 8T bit-cell without assist, and (c) 8T bit-cell with write-assist. VDD=250mV. White: erroneous bit; Black: correct bit.	117

4-4	For a 65nm technology, 10K MC simulation shows that 200mV higher V_{WWL} improves write failure σ/μ by $3.8\times$ at 400mV.	118
4-5	8T bit-cell based design with write-assist enables operation from 1.2V down to 0.37V.	119
4-6	Two options for partitioning the voltage domains in dual supply write-assist technique (a) Level shifters between boundary and logic, (b) Level shifters internal to SRAM.	120
4-7	The schematics of the DCVSL type level shifter used in this design under level shifting operation.	120
4-8	Organization of the 128 kb SRAM. (a) The two-stage sensing of the 32kb block. (b) The schematics of the kth sub-block.	122
4-9	Two stage sensing signals during two back-to-back read operations.	124
4-10	First level sensing (FLS) inverters use a 100 mV VTC shift to optimize time improvement vs. area overhead.	125
4-11	Body-biased SA (BBSA) (a) schematics, and (b) layout. (Not drawn to scale).	126
4-12	BBSA operation explained in detail - Reset phase.	127
4-13	BBSA operation explained in detail - Determine offset phase.	127
4-14	BBSA operation explained in detail - Assign body voltage phase.	128
4-15	Measured BBSA input referred offset voltages before and after calibration.	129
4-16	Improvement in performance by using two stage sensing and BBSA compared to the conventional design.	130
4-17	Measured energy per operation results vs. V_{DD} for the 128kb SRAM designed in 65nm CMOS.	131
4-18	SEM picture of the $0.152\mu^2$ 6T bit-cell. This bit-cell is used in the 0.5Mb SRAM test-chip [7]. <i>Courtesy of STMicroelectronics</i>	132
4-19	FD-SOI transistor cross-section [8].	133
4-20	Flip-well (FW) vs. conventional-well (CW) transistor structures of FD-SOI.[7]	133

4-21	FW vs. CW body connections during no body-biasing.	134
4-22	Write margin μ / σ vs. V_{DD}	136
4-23	Read SNM μ / σ vs. V_{DD}	136
4-24	WVBL under body-biasing under worst case operating conditions. . .	137
4-25	body-bias required for voltage scaling.	138
4-26	Read SNM under body-bias ($V_{DD}=500\text{mV}$).	139
4-27	body-biasing starts after BL discharges.	139
4-28	(a) Schematics, and (b) layout of the 6T bit-cell in the FW FD-SOI process. Layout is not drawn to scale.	141
4-29	Stepwise adiabatic charging decreases energy dissipation by N.	142
4-30	The memory organization of the 0.5Mb SRAM design.	143
4-31	The schematics of the local sense circuit.	144
4-32	Important signals of the SRAM during two read operations.	145
4-33	The schematics of the body control circuit.	146
4-34	The schematics of the replica circuit.	147
4-35	The signals used for stepwise body charging.	148
4-36	Block diagram of the BIST circuit for self-testing.	149
4-37	The simulation setup for the dynamic read stability margin calculation of the 6T bit-cell.	150
4-38	Dynamic read margin simulation results with NOR=8,32 and 128. . .	151
4-39	Die photo for the 28nm FD-SOI SRAM.	152
4-40	Measured energy consumption vs V_{DD} with and without write-assist for the 0.5Mb SRAM test-chip designed in 28nm FD-SOI technology. .	153
4-41	Read to write ratios for different benchmarks.	154
4-42	Shmoo plot for the 0.5Mb SRAM test-chip designed in 28nm FD-SOI technology.	155
A-1	FBR vs. the memory size for 95% yield	169
D-1	A typical buck converter.	177

D-2 Operation of a buck converter in (a) PWM mode control, and (b) in
PFM mode control. 178

List of Tables

1.1	Relevant self-aware system examples.	31
2.1	Test-chip specifications for the 128kb SRAM with an embedded energy monitoring circuit.	68
2.2	Comparison of the proposed energy sensor with recent work.	70
3.1	Test-chip specifications for the self-aware processor SoC designed in 0.18 μ m CMOS.	101
3.2	Comparison of the self-aware test-chip to previous work.	109
4.1	Comparison of the sense amplifier with offset compensation using body-biasing (BBSA) to recent work.	129
4.2	Test-chip Specifications.	152
4.3	Comparison of results with recent publications. *BL-tracked NBL, **suppressed coupling signal for NBL, ***write-recovery-enhancement lower-cell- V_{DD}	157
B.1	Various Bit-cell Assist Techniques for Enhancing Write-ability	171
B.2	Various Bit-cell Assist Techniques for Enhancing Read-ability	172
C.1	Vanilla LEON3 processor core features.	173

Acronyms

6T six-transistor

8T eight-transistor

ADC analog to digital converter

ADDR address input

AHB AMBA Advanced High-speed Bus

BBSA body-biased sense amplifier

BEN body enable

BIST built-in self test

BL bit-line

BLB complementary bit-line

BOX buried oxide

CW conventional well

DAC digital to analog converter

DCVSL differential cascode voltage switch logic

DEC double error correcting

DI input data

DIBL drain induced barrier lowering

DRV data retention voltage

DSU debug support unit

DVFS dynamic voltage and frequency scaling

ECC Error Correction Coding

ESL equivalent series inductor

ESR equivalent series resistor

EOP energy-per-operation

FBB forward body-bias

FBL first-level bit-line

FBR Fail-bit ratio

FD-SOI Fully Depleted Silicon On Insulator

FLS First-level sensing

FPU floating-point unit

FW flip well

GBL global bit-line

GRLIB Gaisler Research IP library

GSM Global System for Mobile

IC integrated circuits

IPS instruction per second

LER line edge roughness

MC Monte Carlo

MCU multi-cell upsets

MLCC multi-layer ceramic capacitor

NBL Negative BL

NOR number of rows

ODA *observe-decide-act*

OOE out-of-order execution

PCU power-control-unit

PD pull-down transistors of the bit-cell

PFM Pulse Frequency Modulation

PG pass-gate transistors of the bit-cell

PU pull-up transistors of the bit-cell

PWM Pulse Width Modulation

RBB reverse body-bias

RBL read bit-line

RDF random dopant fluctuation

RISC Reduced Instruction Set Computer

RSNM Read Static Noise Margin

RW read/write input

RWL read word-line

SA sense amplifier

SBL second-level bit-line

SCE short channel effect

SEC single error correcting

SEU single event upset

SLS second-level sensing

SNM Static Noise Margin

SoC system-on-chip

SPARC Scalable Processor ARChitecture

SRAM Static Random Access Memories

WL word-line

WLEN word-line enable

WWL write word-line

WLUD WL under drive

WVBL bit-line write margin

UTBB ultra-thin body and buried oxide

pchg precharge signal

cSel column select signal

i – cache instruction-cache

d – cache data-cache

C_{STO} storage capacitor

V_{MIN} minimum supply voltage

V_{ofs} comparator offset

I_{READ} bit-cell read current

V_{STO} voltage across storage capacitor

V_T threshold voltage

Chapter 1

Introduction

The semiconductor industry has been faithfully following Moore's law of scaling [9] for almost half a century and the interest for integrating more functionality on a single chip will continue unabated into the foreseeable future [3]. The doubling of transistors in a single die has enabled smaller and faster transistors. This integration enhanced the processing capabilities and features of many handheld devices such as mobile multimedia, wireless sensor nodes and biomedical implants.

This aggressive scaling of CMOS technology has created huge design challenges. Firstly, most of the portable devices are powered up by a battery and the physical limits of electro-chemistry have prevented battery technologies to advance at the same rapid rate as the shrinking of transistor sizes or integrating more transistors on a single die [10]. Thus, the high computational demand depletes the batteries very quickly and energy efficiency is one of the major challenges of today's systems. Figure 1-1 shows the advances of the Li-ion battery energy density over the years. The energy density has increased by $2.85\times$ from 1991 to 2013 whereas the number of transistors on a single die increased by more than 2 orders of magnitude [11, 12].

The ability to harvest ambient energy through energy scavenging is a possible solution for improving battery life-times [13]. However, most practical scavengers can provide up to a few hundreds of μWs and it is still not in the range of providing battery-less solution for many portable systems.

Another challenge of today's systems is the fact that with more integration, they

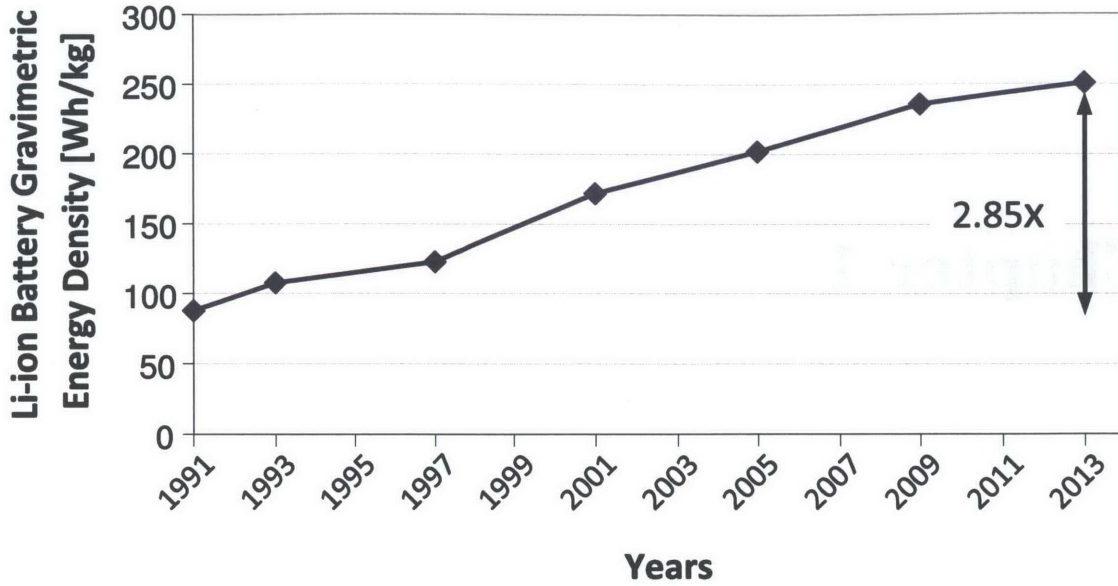


Figure 1-1: Li-ion battery energy density over the years.

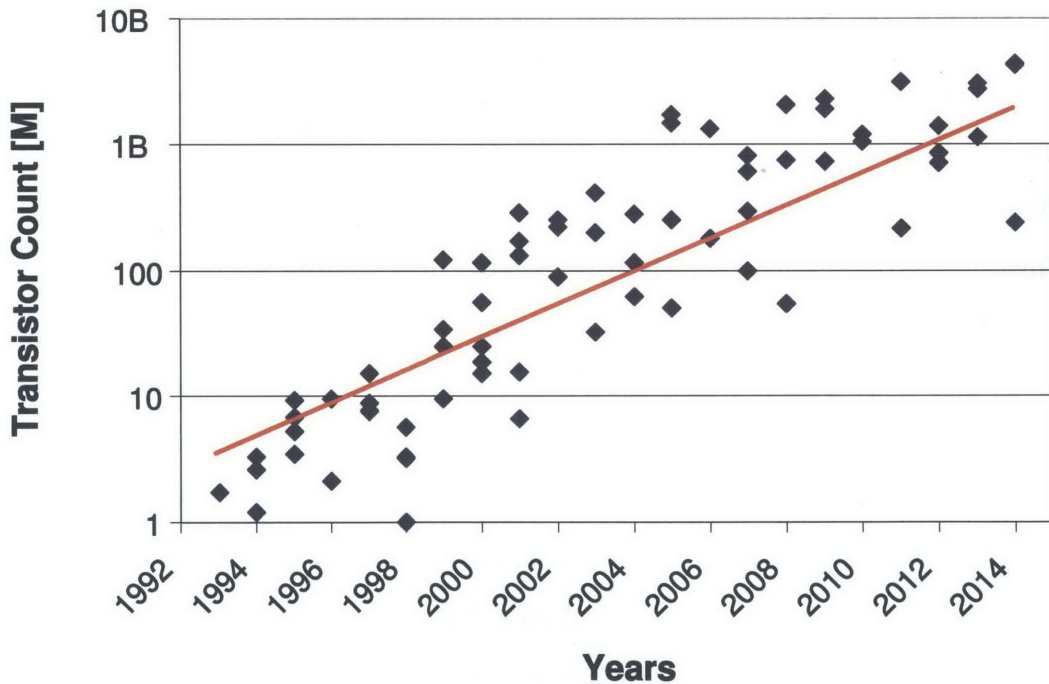


Figure 1-2: Chip complexity over the last two decades.

are getting increasingly more complex [14, 15, 16]. The complexity chart in Figure 1-2 shows the trend in transistor integration on a single chip over the past two decades. Today, it is common to see more than 1B of transistors within a single die. Additional difficulty stems from the fact that those complex systems have to work optimally

under dynamic operating conditions such as temperature and voltage fluctuations, process variations, aging effects, and changing application loads. The complexity and the dynamic nature of the systems create a multi-dimensional design space where it is hard to predict the behavior of the system during the design phase.

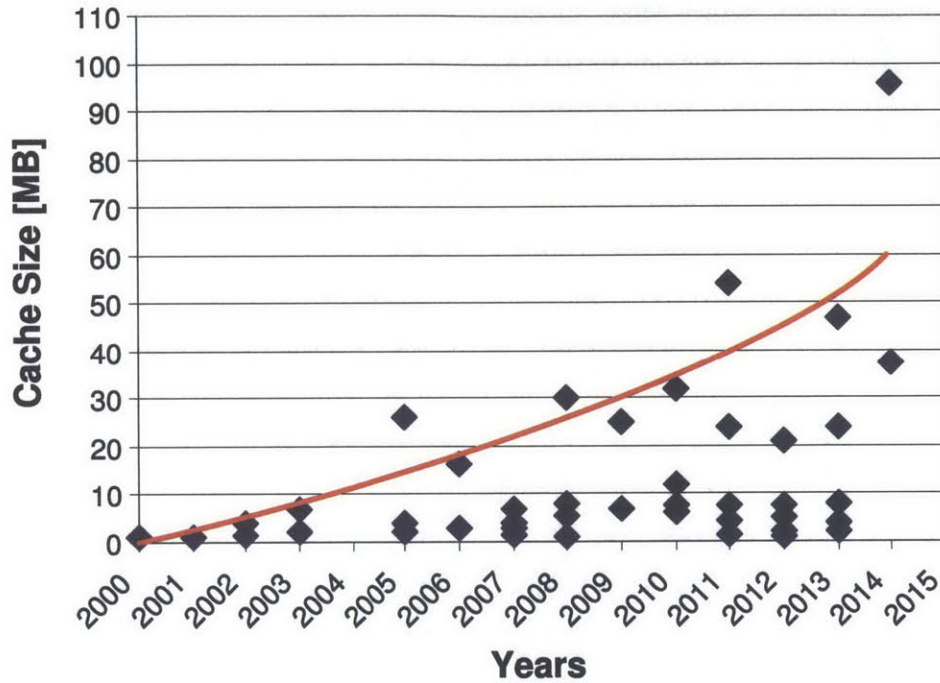


Figure 1-3: Trend of cache integration over the last two decades.

Static Random Access Memories (SRAM) are one of the fundamental building blocks of today’s complex systems. Due to their very regular structure, SRAMs can be developed using aggressive layout rules that address sub-resolution fabrication limits. This way, they provide a high density which is vital for transistor scaling. Secondly, they have a relatively low activity factor since only one word of the memory is being accessed every clock cycle. As a result of those advantages, embedded SRAM size continuously increases over the years.

The size of on-die caches for microprocessor applications over the years can be seen in Figure 1-3. In today’s microprocessors, it is common to see more than 50 MB of on-chip cache design. Recently the Power8 processor has announced to have 96MB on-chip L3 cache [2]. Therefore, embedded SRAMs account for a large portion of the total energy consumption, area and access time of today’s systems and they

will continue to be one of the fundamental building blocks of today’s microprocessor systems.

This thesis provides solutions to complexity and energy-efficiency challenges in a collaboration of circuit and architecture techniques. To enable better architectural decisions, this thesis shows that an energy monitoring circuit which can generate absolute energy-per-operation (EOP) numbers can be beneficial. This circuit idea is extended to a self-aware microprocessor architecture that uses multiple energy monitoring circuits for better energy and performance optimizations. This requires a coherent behavior of circuits, architectures and software. Moreover, as an integral part of today’s system-on-chip (SoC) s, low-voltage SRAM designs are investigated. Particularly, we look into assist techniques to make the SRAMs work at lower voltages.

This chapter begins by briefly explaining self-aware systems and related work in Section 1.1. Then, Section 1.2 talks about dynamic voltage and frequency scaling (DVFS) which is an effective way to reduce total energy consumption. In Section 1.3 challenges of low-voltage SRAM design and previous work are provided. Lastly, Section 1.4 provides thesis contributions.

1.1 Self-Aware Systems

The field of modern electronics sees a continuing trend towards increasingly complex [2, 24, 25] and parallel architectures [26, 27]. Those complex systems generally try to target competing design challenges such as maximizing performance while minimizing power consumption. Furthermore they have to continue working efficiently under changing system dynamics and application loads. These developments require new approaches to design and operate systems that are capable of dealing with complexity and changing behavior. Self-awareness is a growing field of research in computing systems and circuits. Self-aware systems and architectures collect and maintain information about the current state and environment, process the information, decide about the next behavior and adapt themselves if necessary [28, 29].

Today’s systems propose adaptivity in both hardware [18, 19, 17, 21] and in soft-

Table 1.1: Relevant self-aware system examples.

Work	Design Idea	Chip?
Choi, ISCA, 2006, [17]	An approach for processor thread distribution to optimize performance.	No
Albonesi, Computer, 2003 [18]	Adaptive processing for saving energy consumption.	No
Dubach, Micro, 2010 [19]	A control mechanism based on predictive model for adaptivity control.	No
Hoffmann, DAC, 2012 [20]	A self-aware computing model to meet dynamic design goals.	No
Ramadass, ISSCC, 2007 [21]	An energy minimization loop using relative energy measurements.	Yes
Yuffe, ISSCC, 2011 [22]	A power/performance optimized GPU based on model-based observation to adapt cores, L3 cache, DVFS.	Yes
Damaraju, ISSCC, 2012 [23]	A power/performance optimized GPU based on model-based observation to adapt cores, L3 cache, DVFS.	Yes

ware [30] and Table 1.1 summarizes some relevant previous work. One limitation of most of the existing self-aware systems is that they are *closed* [20]. In other words, self-awareness is not accessible by different components of the system. For instance, many hardware-based approaches assume a fixed set of adaptations, which exist exclusively in hardware and are unable to incorporate application specific goals. Similarly, many software based systems assume the hardware is fixed. However, in order to achieve the optimum point of operation, both hardware and software resources need to be available to the decision engine and it should be able to make decisions based on the current information and application targets as a whole.

One example of a close adaptive hardware system is the self-aware system given in [31]. This system does an excellent job of optimizing the total system *throughput* by allocating resources. However, it would not be possible for this system to optimize itself for different design goals such as achieving a desired performance for a specific

application. Similarly, it would not be able to incorporate additional adaptations specified at software level. A second example is the system given in [21]. This system uses a minimum energy tracking loop to coordinate resources to meet the *minimum energy point*. However, this system cannot incorporate additional adaptations such as minimizing the energy consumption while meeting a certain performance requirement.

To avoid sub-optimal conditions of closed adaptive systems, self-aware systems that can coordinate system components are necessary. This requires an interface between hardware and software. An example software computational model is conceptually shown in [32]. This computational model generates a general interface for allowing adaptations that are supported by different system components. However, hardware systems which can support those interfaces are necessary. Furthermore, the hardware system need to be able to continuously monitor its conditions and available resources and it needs to be capable of determining how best to use resources to meet goals given the current system conditions.

In order to be able to change system conditions, it is essential to measure the value to be changed. Some of the important metrics to monitor are performance and energy (or power) consumption. To enable on-fly power and performance optimizations, recent systems leverage power management engines that use monitoring circuits based on energy models [22, 23]. However, different dynamic conditions such as temperature and voltage fluctuations, process variations and aging effects make this a multi-dimensional problem where models cannot fully represent the actual profile. Therefore, energy monitors that can measure the actual energy consumption of the systems are necessary.

The energy monitors have to be able to provide accurate results and should be non-intrusive to circuit operation. Furthermore, they need to require a small area and energy consumption for them to result into net energy savings for a negligible area overhead. In addition to the model based energy sensors, circuits that can measure the actual energy consumption are also investigated in the literature. One example is the energy monitoring circuit which is presented in [33]. However, this circuit only works with a PFM mode DC-DC converter and achieves 20% accuracy only after a

calibration process.

1.2 Energy-efficiency and voltage-scaling

In the evolution of modern electronic devices, energy-efficiency has become one of the bottlenecks of the battery operated integrated circuits such as smartphones. In late 1990s, a Global System for Mobile (GSM) phone contained a simple Reduced Instruction Set Computer (RISC) processor running at 26MHz, supporting a primitive user interface. However, following the trend of recent smart phones and laptops, recent microprocessor systems have become much more advanced and they at much faster frequencies. Battery capacity and the thermal limits imply a power budget of roughly 3W for a smartphone and the available peak power budget for digital circuits is around 2W [3]. This limited power budget requires energy-efficiency in designing the fundamental building blocks of the system. In addition to power management strategies, voltage and frequency scaling is a frequently used method to lower the energy consumption of the digital circuits.

The total energy per operation of a digital circuit can be split into two components: dynamic energy and leakage energy. The dynamic energy component scales quadratically with V_{DD} . On the other hand, the leakage energy component, which is the leakage power that integrates over the time period of operation, increases due to the fact that the circuits are getting slower as the voltage goes down. While it is negligible at higher voltages, the leakage energy component (consisting of gate, junction and subthreshold leakage components) increases exponentially as V_{DD} is decreased close to the threshold voltage (V_T). These opposing trends of active and leakage energy components gives rise to a minimum in the total consumed energy. The minimum energy point lies in the sub-threshold region ($V_{DD} < V_T$) for many systems where devices are operated in weak inversion [34, 35, 36]. The decrease in V_{DD} results into degraded circuit performance. However, this is mostly tolerable since there are large periods of time when the workload required for the digital circuits are much smaller than the peak. Therefore systems generally prefer to trade-off energy

vs. performance.

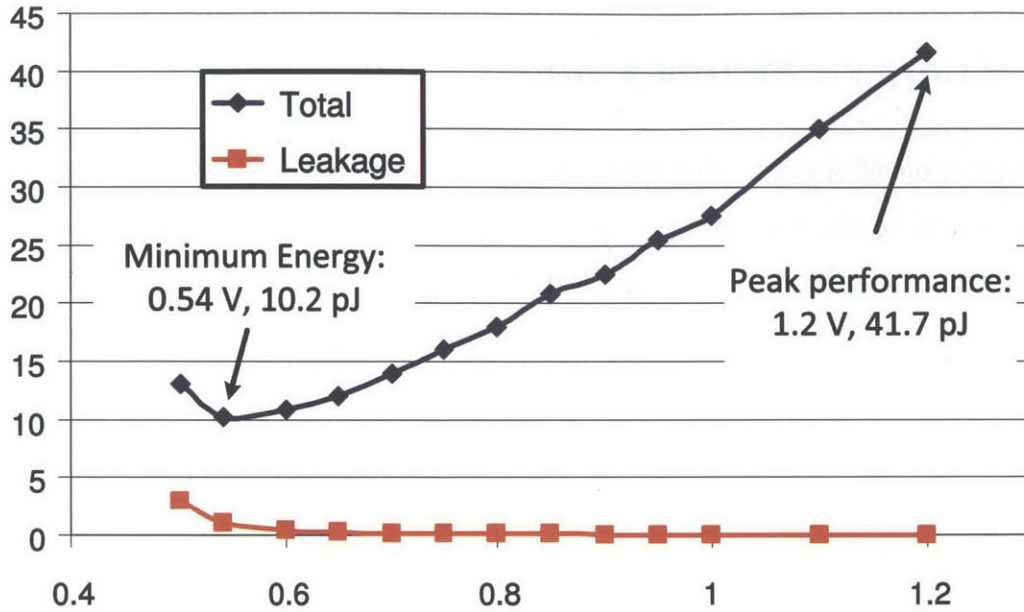


Figure 1-4: Energy and frequency scaling with voltage for a microprocessor application.

Figure 1-4 shows the active and leakage energy profiles of a voltage-scalable microprocessor SoC designed in 65nm CMOS [37]. At the highest voltage supported by the process (1.2 V) the SoC operates at 82.5MHz and consumes 41.7pJ/cycle. At the low-voltage design point of 0.6 V, the energy per cycle is 10.8pJ/cycle at 1.65MHz. The minimal energy per cycle (10.2pJ/cycle) is actually achieved at a slightly lower voltage of 0.54 V. Beyond this point, performance degradation causes leakage power to be integrated over very long access periods and makes leakage energy dominant. Leakage, performance and energy/access plots show that operating at around $V_{DD}=0.5$ V can provide reasonable performance and significant energy savings.

Voltage scaling also has the advantage of lowering the leakage power consumption due to drain induced barrier lowering (DIBL) effect [38]. Especially at deeply-scaled CMOS technologies, where leakage power is significant, voltage scaling can provide exponential savings in leakage power. For instance, at 65nm, leakage power can be reduced by more than 30 \times over the voltage range of 1.2 to 0.5V [39]. Especially for applications such as SRAMs where leakage power is important, voltage scaling can

provide significant energy savings.

1.3 Low-Voltage SRAM Design

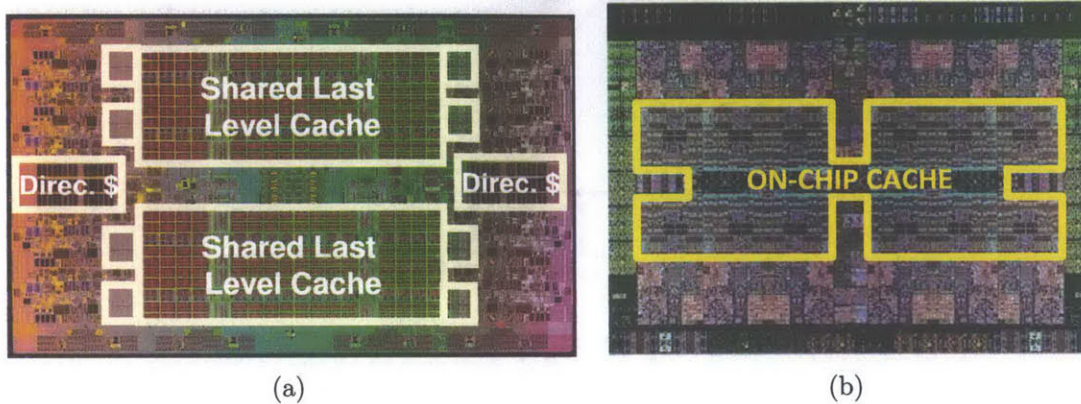


Figure 1-5: (a) Itanium microprocessor with 54MB [1], and (b) Power8 microprocessor with 96MB [2] on-chip caches.

In today's SoCs and microprocessors, embedded SRAMs comprise a large portion of total chip area and energy consumption. Figure 1-5 shows two examples of modern microprocessors where embedded SRAM caches dominate the total chip area. Itanium processor [1] has a total of 54MB on chip cache and 50MB of this is SRAM memories. The Power8 microprocessor [2] has a record of 96MB on-chip L3 caches.

To increase memory density, memory bit-cells use the smallest devices in a technology. Figure 1-6 shows SRAM bit-cell area scaling over the last two decades and from the figure it can be seen that SRAM bit-cells follow the trend of Moore's law of scaling.

In advanced CMOS nodes, the predominant yield loss comes from the increase in process variations. Some examples are the local variations due to random dopant fluctuation (RDF) and line edge roughness (LER). Those process variations degrade performance and increases leakage in logic circuits. However, their impact on SRAM is much stronger. First of all, the fact that SRAM bit-cells are very small in terms of area, makes SRAMs more vulnerable to variations [40]. Secondly, these effects degrade SRAM functionality as the supply voltage is reduced. Figure 1-7 shows that

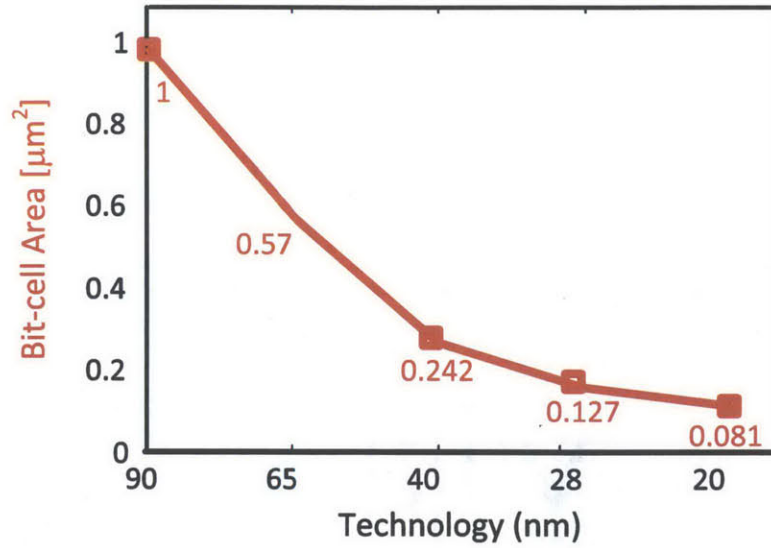


Figure 1-6: Bit-cell area scaling over technology nodes. [3]

V_T variation of SRAM devices increases significantly with scaling, which poses a major challenge for SRAM design [4].

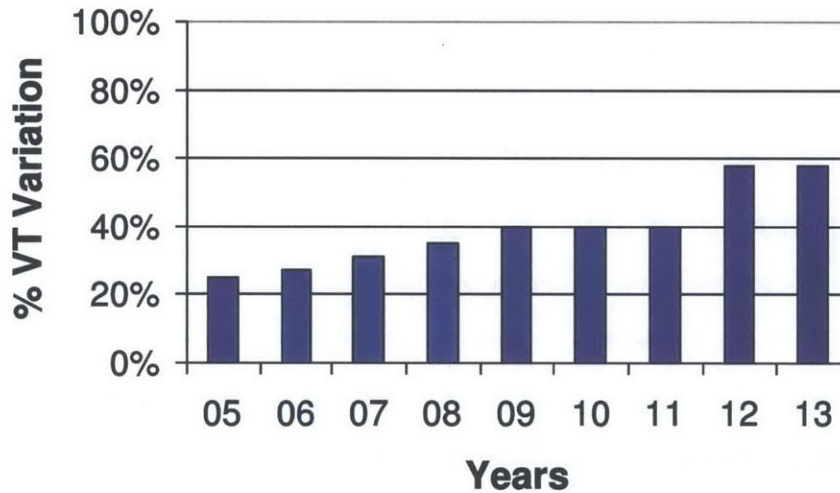


Figure 1-7: SRAM devices V_T variation scaling trend [4].

1.3.1 Sources of Variation

Variation is the deviation from intended values of structure. The electrical performance of modern integrated circuits (IC) is subject to different sources of variation. For the purposes of circuit design, the sources of variation can be categorized in two

classes [41, 4, 42]:

- Global Variation: This is also known as die to die or inter-die variation. This type of variation affects all devices on the same chip in the same way. For example all the transistor gates can be shorter than the nominal value.
- Local Variation: This is also known as within-die or intra-die variation. This type of variation corresponds to variability within a single chip, and may affect different devices differently on the same chip. For example devices in close proximity may have different V_T than the rest of the devices.

1.3.2 SRAM Failure Mechanisms

SRAMs experience three types of failures:

- Hard failures: The hard failures are due to catastrophic defects which can cause permanent damage to the bit-cell and they effect the product yield significantly (Appendix A). Some hard failure examples are metal bridges or missing vias. Redundancy in rows, columns or banks improves hard-error susceptibility [43].
- Soft failures: SRAMs are susceptible to single event upset (SEU) s which arise due to radiation particles that hit the silicon substrate. With technology scaling, transistor sizes are shrinking and SEU rate per bit-cell is decreasing. However, the memory capacities are increasing which still makes soft-errors important for SRAM design. In literature, different techniques are used to cope with soft errors. Some examples are using SOI technology, column interleaving and using Error Correction Coding (ECC) [44].
- Stability failures: This is the main limiting factor for SRAM supply voltage scaling, or minimum supply voltage (V_{MIN}) [45, 46]. Therefore, these failures are the main focus of this thesis. There are four types of stability failures:
 1. Read access failures.
 2. Read upsets (or read stability failures).

3. Write failures.
4. Hold (or retention) failures.

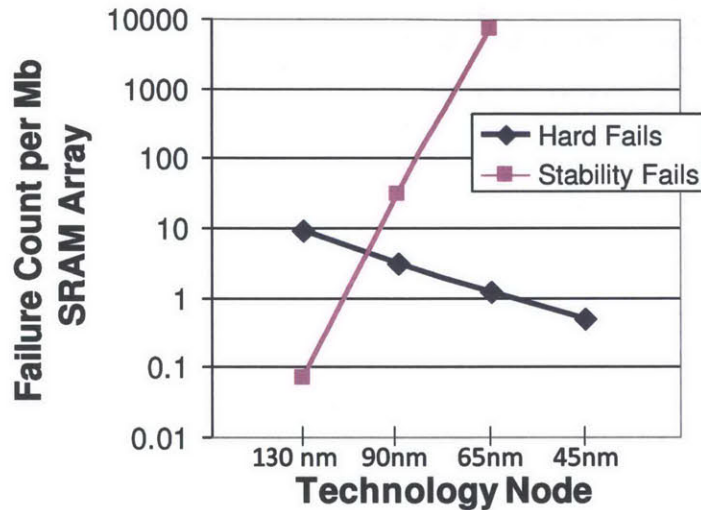


Figure 1-8: SRAM failures due to hard errors and stability issues versus technology nodes [4].

Figure 1-8 shows the effect of technology scaling on SRAM failure probability for the technology nodes between 130nm down to 45nm. Hard failures decrease due to the reduction in bit-cell area and improvement in defect density. However, as bit-cell size is reduced by $2\times$ at every technology node, process variations increase significantly and become the dominant cause of bit-cell failures [40].

The conventional six-transistor (6T) bit-cell can be seen in Figure 1-9. M1 and M2 are the pass-gate transistors of the bit-cell (PG), M3 and M4 are pull-down transistors of the bit-cell (PD), and M5 and M6 are pull-up transistors of the bit-cell (PU). For correct operation, the 6T bit-cell needs to ensure that the contents of the cell are not altered during a read operation and the cell should be able to quickly change its state during a write operation. These two operations require conflicting trade-offs for read and write which makes it hard for the 6T bit-cell to provide stable read and write operations [47, 4].

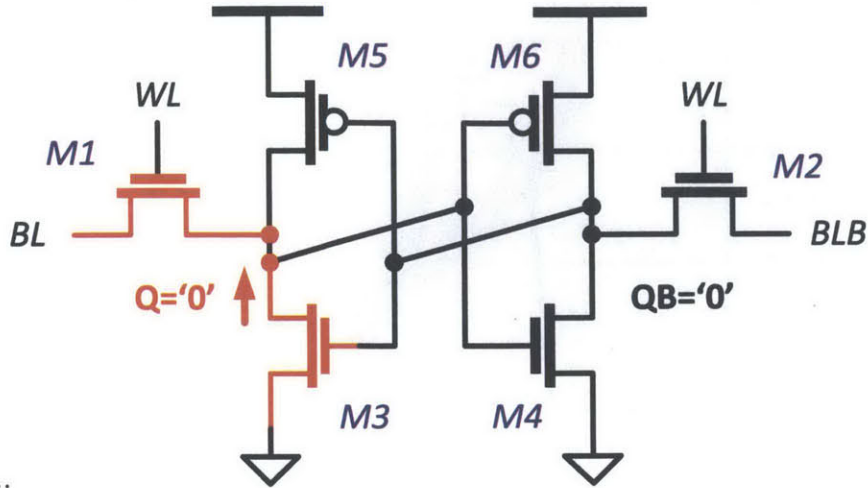


Figure 1-9: 6T bit-cell in read operation.

Read Stability Failure

During read operation of the 6T bit-cell, bit-line (BL) and complementary bit-line (BLB) nodes are precharged to V_{DD} and then, word-line (WL) is asserted. Due to the voltage divider between M1 and M3, the storage node of Q slightly increases. If Q rises close to the V_T of M4, the cell might flip its data. Therefore for stable read operation, PD needs to be stronger than the PG.

Static Noise Margin (SNM) is used to quantify the robustness of the bit-cell against stability failures [48]. Read Static Noise Margin (RSNM) is calculated by finding the largest square which fits inside the voltage transfer characteristics (VTC) of the two back to back inverters which is illustrated in Figure 1-10. The read stability failure occurs if $SNM \leq 0$. This error can occur anytime WL is enabled even if the bit-cell is not accessed for read or write. This situation happens for half-selected bit-cells which experience a dummy read operation.

Dealing with read stability failures is one of the biggest challenges of SRAM design and it has been extensively studied in the literature. Important assist techniques can be summarized as WL under drive (WLUD), using a lower BL capacitance, precharging bit-lines to a lower voltage, body-biasing and read and write back. The WLUD technique which is one of the most commonly used read assist technique and it has been used by [49, 50, 51, 52, 46]. In this method, the WL voltage is lowered

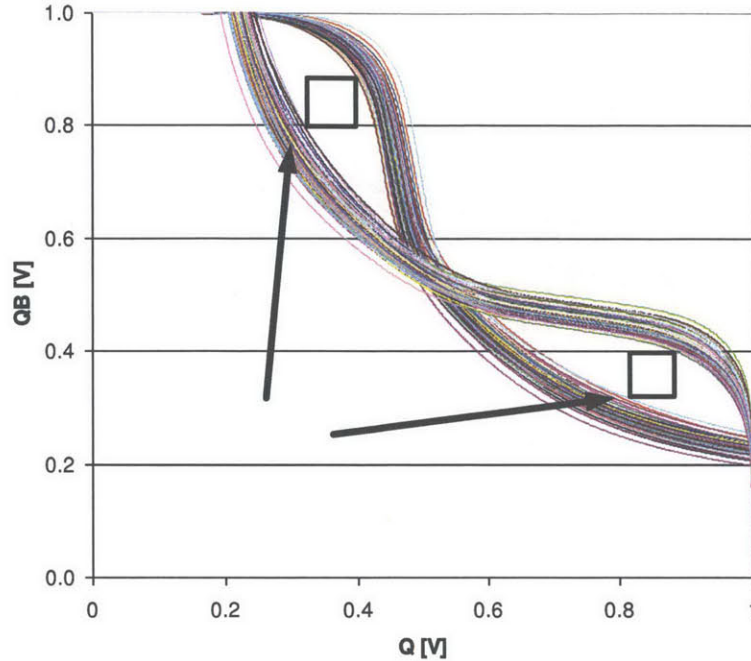


Figure 1-10: RSNM butterfly curves for a 28nm bit-cell for 100 MC runs.

during read operation to decrease the strength of the PG transistor. Another method is using a lower voltage to precharge the BL and BLB nodes [53, 54, 55]. Furthermore, a third technique is to take advantage of the dynamic nature of read operation and using reduced bit-line capacitance to reduce read disturb [56, 57].

Read Access Failure

A typical read path of a 6T bit-cell based SRAM is given in Figure 1-11 and it is the critical path in SRAM memories [58]. This example shows a memory structure with 2 to 1 column interleaving where one sense amplifier (SA) is shared between two columns and either one of the columns is selected during a write or a read operation.

As explained in the previous section, at the beginning of the read operation BL and BLB are precharged to V_{DD} . Then, the active column is selected depending on the address input. Then, the WL is asserted. This WL activation is a small period of time which is determined by the BL and BLB capacitance and the bit-cell read current (I_{READ}). After WL is asserted, a voltage differential (V_{diff}) starts to build up between BL and BLB by I_{READ} depending on the stored data inside the bit-cell. The

V_{diff} has to be greater than the SA input offset voltage for correct read operation. Read access failure occurs if I_{READ} drops below the designed limit and the bit-cell cannot create the necessary V_{diff} . This can happen if one or both of the PG and PD transistors become weaker due to variation. Secondly, it can also happen if the SA input offset is high. This type of failure impacts the memory performance since the WL activation time is around 30% of the memory access time [59].

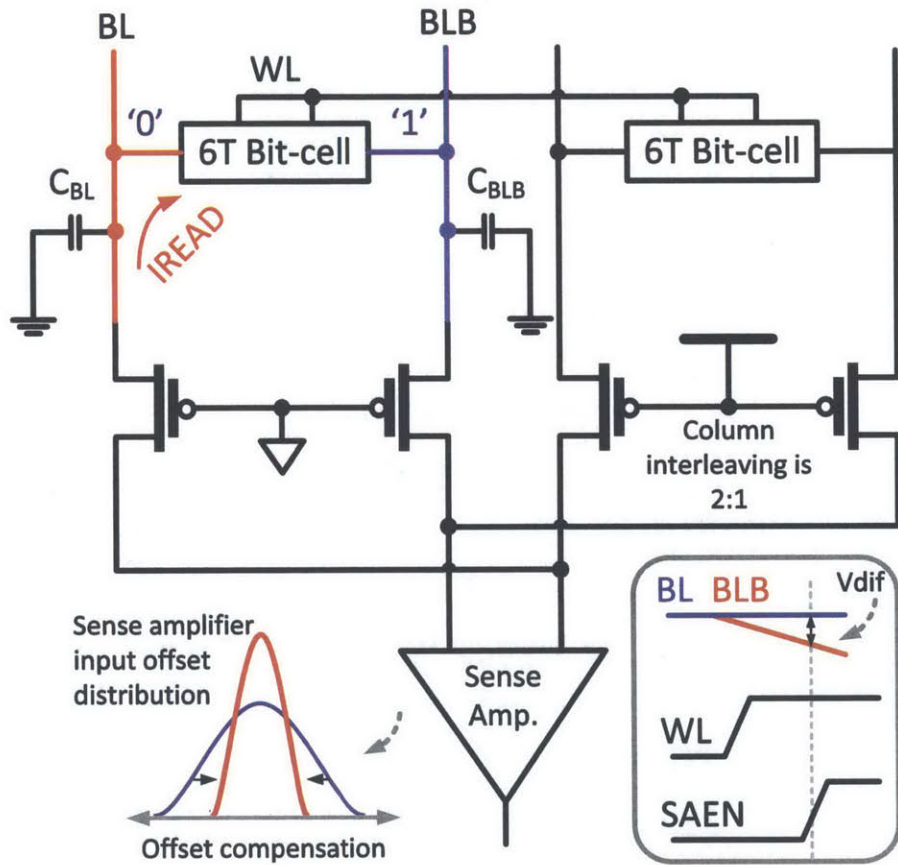


Figure 1-11: A typical SRAM read path.

Offset compensation for the SAs is a technique to tighten the SA input offset distribution. This helps with both performance and energy consumption due to BL discharge. There has been extensive research on SA offset reduction in the literature. The authors in [60, 61] used tunable pseudo-differential SAs which select the reference voltage among different voltage differentials to compensate for offset. Authors in [62] use SA redundancy and select better performing SA during the calibration process.

Write Stability Failure

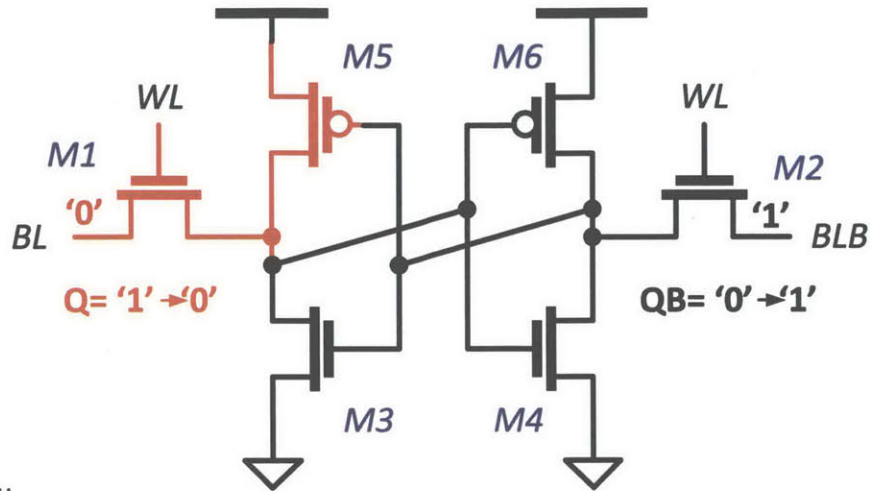


Figure 1-12: 6T bit-cell in write operation.

The ability to write the new data into the bit-cell is referred as the write stability. 6T bit-cell in write operation can be seen in Figure 1-12. At the beginning of the write operation, depending on the data to be written, BL or BLB is pulled down to '0' through write drivers (BL in the case of the figure). Then, the WL is asserted. Thus, M1 turns on and pulls Q node from '1' to '0'. If the voltage of Q falls below $V_{DD} - V_T$ for M6, positive feedback begins for correct write operation. However, before the positive feedback occurs, M5 is turned ON and it is holding Q at '1'. Therefore, for stable write operation, PG needs to be stronger than the PU transistor. Write failure can also happen if the WL pulse is not long enough for the bit-cell to flip the data.

To quantify write stability, one of the most widely used metrics is the WSNM, which uses a similar concept to RSNM. In this approach, VTC of the two sides of the bit-cell are obtained from a DC simulation, while BL and BLB are driven to the write operation conditions. For a successful write operation, there needs to be only one cross-point between the VTCs of the inverters which means that the bit-cell is monostable. Another write stability metric is bit-line write margin (WVBL), which can be calculated using a DC simulation. In this method, bit-cell is configured in write operation with BL connected to V_{DD} and BLB is swept from V_{DD} to ground. As the voltage is swept, the internal nodes flip and the BL voltage at the flip is defined

as the WVBL.

For more realistic stability measurements, dynamic effects need to be considered. For a read operation, performing DC analysis such as SNM simulations create pessimistic results. This is due to the fact that during the read operation, BL voltage discharges, which results into a reduced stress on the storage nodes. Therefore, dynamic read margin simulations are proposed in [63, 57, 64]. On the other hand, the dynamics of the write operation depend on the WL pulse width, which is assumed to be infinite during DC simulations. Therefore, DC analysis create optimistic write margins, and WL pulse width needs to be taken into consideration for a more realistic margin calculation.

To address the challenge of decreased write-margins, several techniques are investigated in the literature. Some of those techniques can be summarized as lowering the bit-cell supply voltage, WL boosting, Negative BL (NBL) and body-biasing. Lowering the bit-cell supply voltage improves write margin, since PMOS PU gets weaker. Several implementations are proposed to either totally collapse or reduce the bit-cell supply voltage during a write cycle [51, 65, 66, 67]. Another write assist technique is the NBL write assist. This helps write operation since it improves the strength of the NMOS PG by applying a small negative voltage to its source [68, 69, 70, 71]. The second important technique is the WL boosting. This technique is typically used with eight-transistor (8T) bit-cells since they do not suffer from column-interleaving problem [37, 72]. Boosting WL voltage improves the strength of the PG transistors and improves write-ability.

Figure 1-13 shows the illustration of various dynamic peripheral assist techniques that are mentioned in this section. The assists are categorized for read and write enhancement separately. Specifically, for read-ability enhancement, we are illustrating techniques such as dynamic elevation of VDD, WLUD, using body-biasing, and precharging the BLs to a lower voltage. Similarly, for write-ability enhancement, we are illustrating dynamic reduction of VDD or elevation of VSS, WL boosting, and NBL.

A detailed summary of various bit-cell assists are listed in Appendix B.

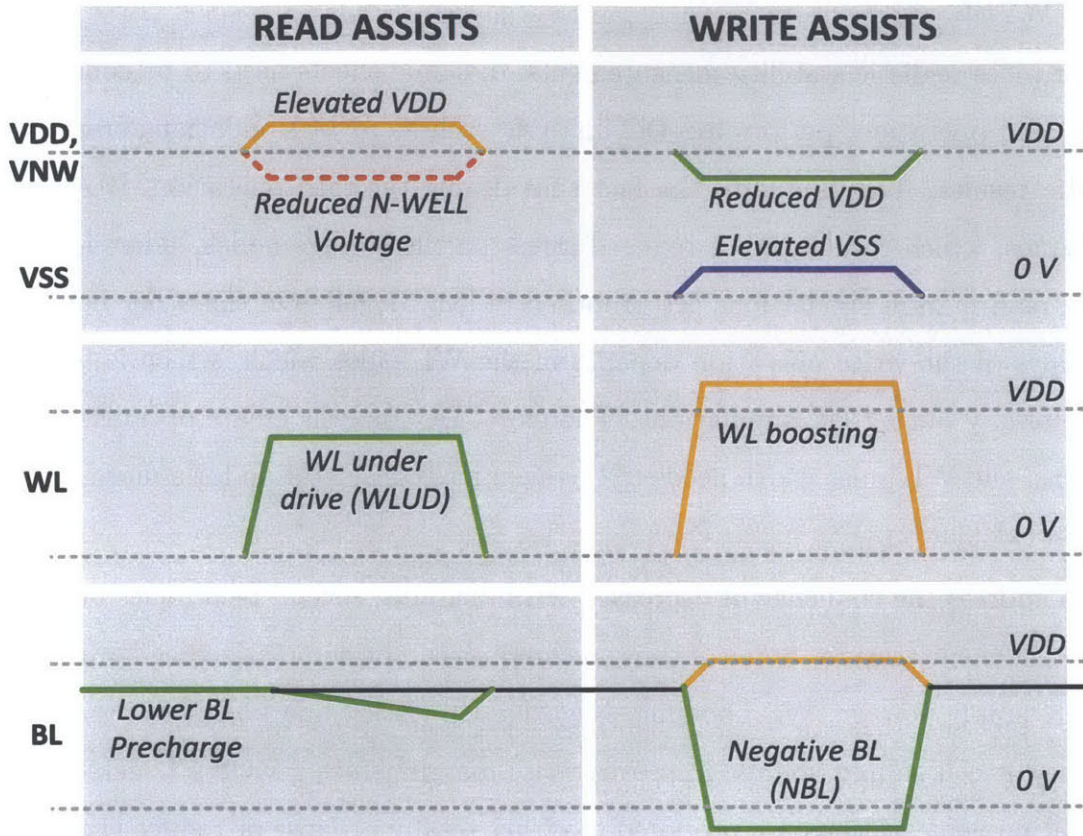


Figure 1-13: Illustration of various peripheral assist techniques. A similar diagram is illustrated in [5].

Data Retention Failure

As explained in Section 1.2, reducing V_{DD} is one of the most popular ways of increasing the energy efficiency of digital circuits. In SRAM circuits, the data retention voltage (DRV) defines the minimum voltage under which the data of the SRAM would not be retained correctly. Due to transistor variations, the cell might be imbalanced and the margin for the retention is measured by using SNM in standby where WL is not asserted. If the cell is imbalanced, its DRV tends to be higher.

1.4 Thesis Contributions

This thesis proposes self-aware system design to enable dynamic power/performance optimizations. Absolute energy numbers are measured using a proposed energy mon-

itoring circuit. The thesis also has a focus on low-voltage SRAM design techniques. It investigates two bit-cell topologies (6T and 8T), and proposes different assist and sensing techniques for them. The main contributions include:

1. **Embedded Energy Monitoring Circuit For SRAM**

For self-awareness, we are suggesting to use the knowledge of runtime energy consumption of the key blocks within the system. In Chapter 2, an energy monitoring circuit is demonstrated for an SRAM application. In this chapter, methodology, design and challenges of the embedded energy monitoring circuit are explained in detail. A test-chip prototype is demonstrated in 65nm LP CMOS [72]. Proposed energy monitor can measure the actual energy consumption rather than estimating it based on a model. Measurement results show that, compared to recent energy (or power) monitors, the proposed circuit achieves the same or better accuracy (10%) without requiring a calibration phase. Moreover, it restricts the voltage drop to a fixed ΔV that can be set by 10mV step sizes. This is important for circuit robustness for many applications including SRAMs. The energy monitoring circuit requires a small area and power overhead and it is non-intrusive to circuit operation. Therefore, it can be used in a system design.

2. **Self-Aware Processor SoC Design Using Multiple Energy Monitors**

This project was implemented by a group of students. I was responsible from many blocks as well as the complete hardware design implementation. Below is a list of contributions:

- I designed the DC-DC converters with embedded energy sensors.
- I was also responsible from the design of the voltage-scalable and reconfigurable SRAMs with Mahmut. However, the cache architecture was done by Eric.
- The system level architecture, LEON3 modifications and simulations required a joint effort, although they were mostly performed by the architecture students (Sabrina, Eric, Jason).

- I did the synthesis of the design, back-end of the test-chip, and the board design.
- I did the initial setup of the board and testing, tested the custom design blocks and performed the initial testing of the processor.
- The software decision engine (SEEC) related work was performed by the architecture students (Hank, Jason, Sabrina).
- Me and Sabrina were responsible from the asynchronous boundary design between core and the custom-design circuits.

Chapter 3 extends the idea of absolute energy monitoring to a self-aware processor design. The system is designed to support the software decision engine, SEEC, which was shown conceptually before [20]. It exposes embedded energy monitoring circuit readings and circuit adaptations to SEEC. This exposure can allow SEEC to coordinate reconfigurable hardware and DVFS actions to meet application specific goals during runtime.

A test-chip prototype is implemented in $0.18\mu\text{m}$ CMOS. The energy monitors are embedded into on-chip DC/DC converters and bring minimal power ($<0.1\%$) and area ($<1\%$) overhead. The custom blocks are designed to be voltage scalable from 1.8V down to 0.6V. Low-voltage SRAM operation is possible thanks to the use of 8T bit-cells and write-assists. The reconfigurable d-cache design enables changing the set-associativity (1-4 sets) and size (1kB-4kB per set) to adapt to compute- versus cache-bound phases of applications. Cache reconfiguration can be performed in < 3 clock cycles including tag invalidation, which takes only one clock cycle. These hardware features enable SEEC to dynamically adapt the microprocessor to meet performance and energy goals. Measurement results show that up to $8.4\times$ energy savings can be achieved with DVFS and self-adaptation [73].

3. Low-Voltage SRAM Design

SRAMs are important building blocks of today's systems and they account for a

significant portion of the total energy consumption and area. Therefore, Chapter 4 analyzes SRAM design especially for low-voltage operation. The trade-offs between 6T and 8T bit-cell topologies are investigated. Assist techniques are proposed for both designs. Moreover, an offset compensated SA is proposed.

Firstly, an 8T bit-cell based SRAM design is demonstrated in a test-chip prototype in 65nm LP CMOS. This SRAM uses WL boosting technique to improve write-ability at low-voltages and achieves voltage scaling down to 0.37 V. However, 8T bit-cell is not compatible with column interleaving. Since SAs cannot be shared with the neighboring column, the SA area is limited. Therefore, we used a new SA offset compensation circuit to provide offset reduction. The proposed circuit utilize body-biasing to reduce input offset by $2\times$.

Although 8T bit-cell promises low-voltage operation, it requires a 40% larger area compared to its 6T counterpart. Therefore, in a second test-chip prototype, we designed a 0.5Mb SRAM using a 28nm FD-SOI technology with industry sized 6T bit-cells. This test-chip demonstrates read and write-assist techniques that leverage extensive body-biasing capability of FD-SOI technology. It is measured to operate down to 0.43 V. This is the lowest voltage achieved for a 6T bit-cell based SRAM design in a 28nm process.

Chapter 2

Energy Sensing Circuit for an SRAM Application

Self-aware systems require observation blocks to monitor the important metrics of the system such as temperature, voltage, battery charge and energy (or power) consumption [20]. Modern systems offer different design options to sense run-time energy or power consumption. One popular way is to predict the value rather than actually measuring it [22, 23, 6].

An example power monitor is given in SandyBridge processor's power-control-unit (PCU). Figure 2-1 shows the block diagram of this processor with its major functional blocks, the PCU, and the interconnect [6]. This PCU requires a combination of complicated state machines, an integrated microcontroller, and thermal sensors to perform power management. This way, it estimates the power consumption based on the thermal readings and voltage information rather than actually measuring it. However, in aggressive CMOS technologies, variation in operating conditions and applications introduce a multi-dimensional design space which is difficult to model correctly. For example, for a 65nm RISC processor, the power can vary by $26.7\times$ among different instructions and extreme corners [74] which motivates measuring the power (or energy) consumption rather than predicting it. Thus, this chapter focuses on the design and challenges of an embedded energy monitoring circuit that can measure the absolute energy-per-operation (EOP) of a circuit.

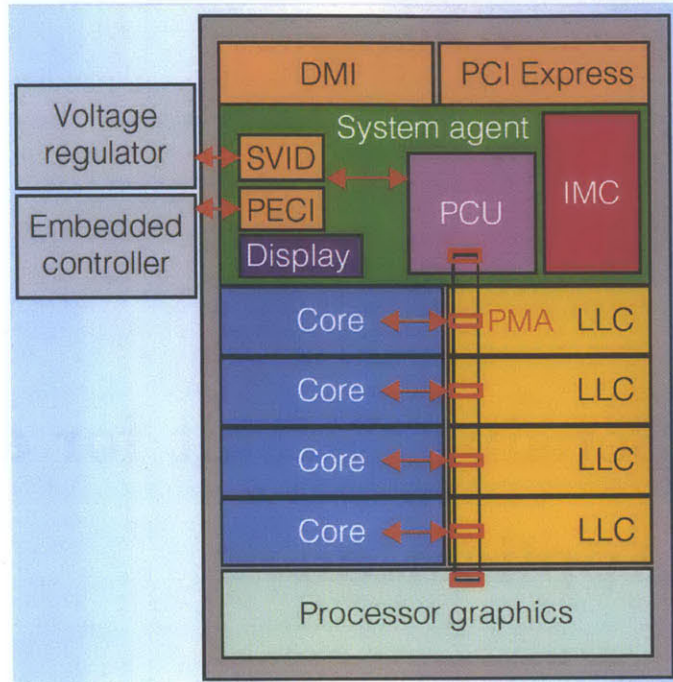


Figure 2-1: SandyBridge power management architecture [6].

As mentioned in Section 1.3, SRAM circuits are one of the fundamental building blocks of modern SoCs and their energy consumption is significant compared to the total energy consumption. Thus, for complex systems with significant size of on-chip caches, the knowledge of absolute EOP of the SRAM blocks can be very useful for system level power and performance optimizations. Therefore, in this chapter, we will demonstrate the idea of energy monitoring for an SRAM application.

First, a prototype test-chip is fabricated using a 65nm LP CMOS process. In this test-chip, the operation of the energy monitoring circuit is demonstrated for a 128 kbit SRAM circuit [72]. The design of the SRAM will be explained in Chapter 4. The idea of energy monitoring is extended to a self-aware processor system with multiple energy monitors in Chapter 3. This system is designed to reconfigure itself based on energy information from multiple embedded energy monitors [73]. This second test-chip prototype is designed in 0.18 μ m CMOS process.

2.1 Overview of Energy Sensing Concept

Different ways have been offered to measure energy consumption. One method is to measure the voltage drop across parasitic resistance along the power delivery path such as the parasitic resistance of the filtering inductor. This method is useful for high power applications of the order of tens of watts. However, to be able to measure load currents around μ watts across the small parasitics, power hungry amplifiers would be needed which is not desirable.

Another way to measure the EOP would be estimating the energy lost across a storage capacitor (C_{STO}) as proposed in [75]. If voltage across storage capacitor (V_{STO}) during normal operation is V_1 and it becomes V_2 after N clock cycles, V_1 and V_2 can be digitized using high-precision analog to digital converter (ADC) s. Then, energy consumed by the circuit under test can be calculated as $C_{STO}(V_1^2 - V_2^2)/2N$ with the help of digital multipliers and dividers. For this method to be practical, the voltage ripple across the storage capacitor due to energy monitoring has to be small. Otherwise, it can degrade the circuit robustness and performance. For a nominal voltage of 1.2V, 10mV can be an acceptable voltage ripple [72]. This would require an ADC with a precision around 1mV. 1mV precision can be detected using at least an 11 bit ADC and it is clearly not an energy efficient solution.

On the other hand, a solution that relies heavily on digital circuitry can bring small overhead in terms of energy. An elegant way of energy sensing is proposed in [75] by observing $V_1^2 - V_2^2$ can be simplified as $2V_1(V_1 - V_2)$ assuming $(V_1 - V_2)$ is small. This sensor uses mostly digital circuits to keep the energy overhead small. However, it does not strictly limit the voltage drop which can result into bit failures for an SRAM application. Also, it measures the relative energy information but not the absolute energy.

In this thesis, an energy monitoring circuit that can measure the absolute EOP during runtime is proposed. The energy monitor does not require a calibration process. The first implementation, which is the focus of this chapter, requires an off-chip C_{STO} ; however, in the next chapter a second version of energy monitors which are

embedded into DC-DC converters will be shown. Those monitors share the off-chip filtering capacitor of the DC-DC converters, thus they don't require an extra off-chip capacitor. The energy monitors are measured to achieve 10% accuracy and they require a small area overhead. Furthermore, they strictly restrict the ΔV drop, which is important for circuit robustness.

A simplified version of the energy sensing concept is illustrated in Figure 2-2 and it can be summarized as follows:

1. As illustrated in phase 0 of the figure, C_{STO} is connected to the power supply through a switch. During the normal operation (when energy sensing cycle is not active) the switch is closed and energy is being delivered from the supply. At this time, the voltage across C_{STO} is equal to V_{DD} and the energy stored across the capacitor is equal to $1/2 \times C_{STO} \times V_{DD}^2$. It should be noted that the IR drop across the switch needs to be kept small. For high-power applications, this can require a large switch whereas, for this application, it did not bring any limitation.
2. When energy sensing cycle starts, the switch turns OFF and the circuit under monitoring (SRAM in this case) is disconnected from the power supply. Therefore, the capacitor starts to power up the SRAMs and the voltage across it starts to drop with a slope depending on I_{SRAM} and the value of C_{STO} . The voltage drops from V_{DD} to $V_{DD} - \Delta V$ in N clock cycles (phase 1). The energy stored across the capacitor becomes $1/2 \times C_{STO} \times (V_{DD} - \Delta V)^2$. Therefore, the total energy consumed by the SRAM can be calculated as $1/2 \times C_{STO} \times [V_{DD}^2 - (V_{DD} - \Delta V)^2]$. Since $a^2 - b^2 = (a - b)(a + b)$, this equation is actually equal to $1/2 \times C_{STO} \times \Delta V \times (2V_{DD} - \Delta V)$.
3. As pointed out, for a robust SRAM operation, ΔV needs to be kept much smaller compared to V_{DD} . Therefore, the following approximation would bring a negligible error [75]: $(2V_{DD} - \Delta V) \approx 2V_{DD}$. This simplifies the energy equation to $C_{STO} \times \Delta V \times V_{DD}$.
4. Since this energy is consumed in N clock cycles, $EOP = C_{STO} \times \Delta V \times V_{DD}/N$.

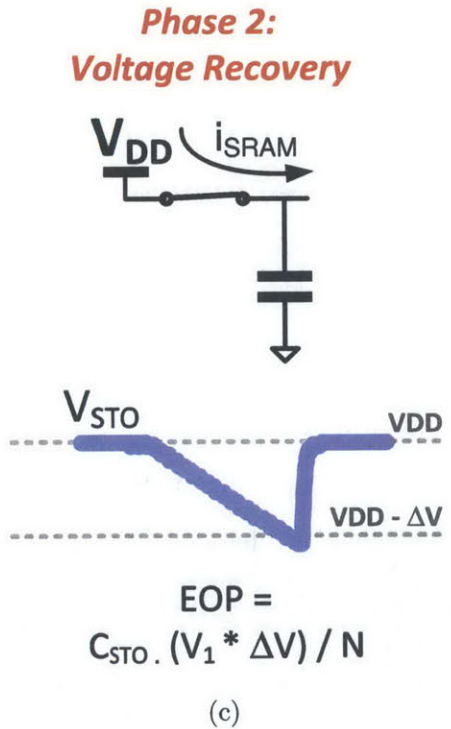
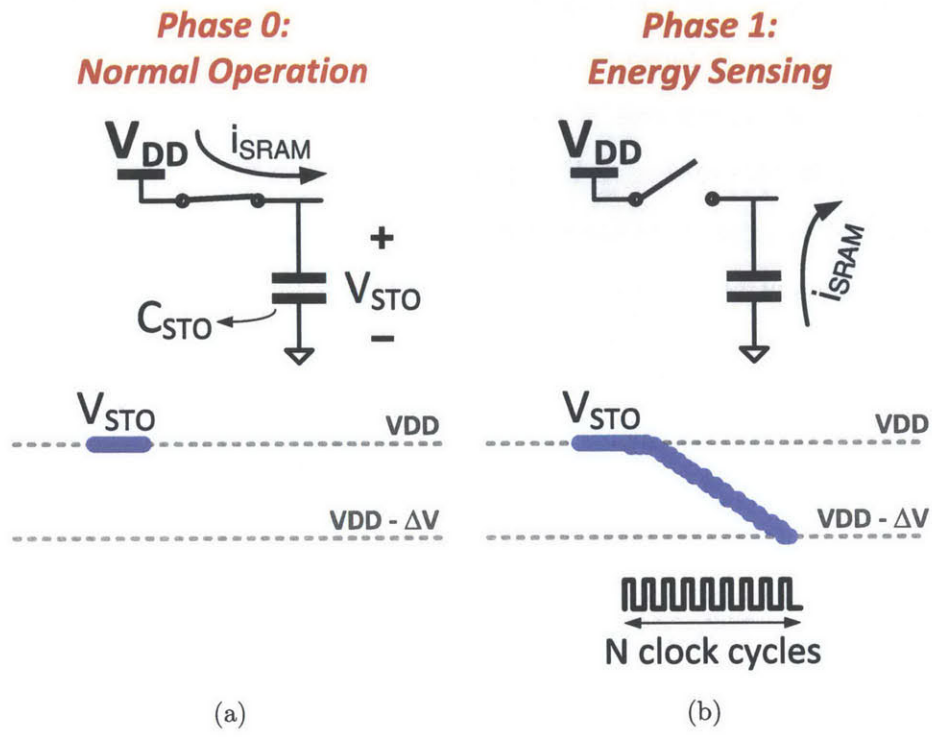


Figure 2-2: (a) Normal operation, (b) Energy sensing cycle, and (c) Voltage recovery.

5. One thing to note here is that in this equation, C_{STO} and V_{DD} are fixed quantities so the two unknowns are the number of clock cycles, N ; and the voltage drop, ΔV . By fixing one of those and observing the other one, EOP can be measured. In this design, ΔV is fixed. It is important to note that, for better accuracy, C_{STO} can be measured in a calibration process to capture its temperature coefficient and parasitics. This is not implemented in this work, although a technique will be suggested.
6. In phase 2, the switch turns back ON and the voltage rises back from $V_{DD} - \Delta V$ to V_{DD} . At the same time, EOP is calculated.

2.2 Energy Monitoring Circuit Challenges

For it to be meaningful in a system application, the energy sensing circuit needs to be almost free. First of all, it should not bring a large energy overhead so that its energy is much smaller compared to the savings it brings. Secondly, it is desirable that it does not require any extra off-chip components or voltages. Thirdly, in order not to complicate the system, the sensor needs to be non-intrusive to the operation of the circuit under test and it is desirable that it does not require any extra calibration process. This way, it can be designed independently. Also, its area needs to be small, especially for applications that can benefit from using multiple energy monitors. Lastly, it needs to provide accurate results.

As pointed out, the energy sensing concept proposed here requires some approximations. Secondly, due to system non-idealities, the accuracy can be degraded. Some of the major challenges of the energy monitoring circuit design as well as important error sources are explained below.

2.2.1 Effect of ΔV Drop on SRAM Operation

The energy sensing concept is based on the supply voltage to decrease across the circuit under test (which is SRAM in this case) during energy monitoring. This

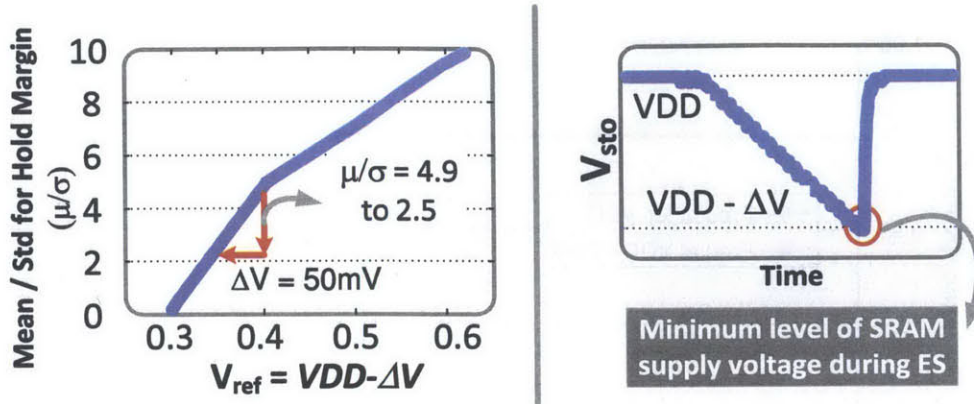


Figure 2-3: Selection of ΔV at low voltages for cell stability.

voltage ripple needs to be non-intrusive to SRAM operation; on the other hand, the SRAM robustness highly depends on the supply voltage. Especially at low voltages, noise margins are very small and even a small voltage drop can result into a severe number of bit-failures as explained in Section 1.3.

Secondly, at low voltages, transistor currents get smaller. This results into degraded performances which can result into timing failures. Therefore, ΔV drop needs to be precisely controlled in order not to degrade SRAM stability and not to create timing failures during energy sensing. The required frequency adjustment during energy sensing is shown in Figure 2-4. At 0.4V, this requires operating the circuit with a 20% lower frequency, during energy sensing.

In Figure 2-3, the μ / σ of the retention SNM distribution can be seen. This ratio is calculated by running 10K Monte Carlo (MC) analyses on the SRAM bit-cell while the PG transistors are kept OFF (Figure 1-9). In this analysis, the SNM distribution is assumed to be Gaussian and μ and σ numbers are calculated accordingly. In this design, μ / σ is kept higher than 4.5 at all times. This value needs to be determined depending on the SRAM size used, how much row and column redundancy is tolerable and the yield target. (Detailed yield calculation for SRAM circuits is explained in Appendix A.) The x-axis is the V_{DD} span from 0.6V down to 0.3V. As it can be observed from the figure, when the V_{DD} is 400 mV and the ΔV is 50 mV, this voltage drop would decrease the μ / σ of the retention SNM from 4.9 to 2.5. This would result into drastic degradation of robustness due to the inability of many bit-cells to

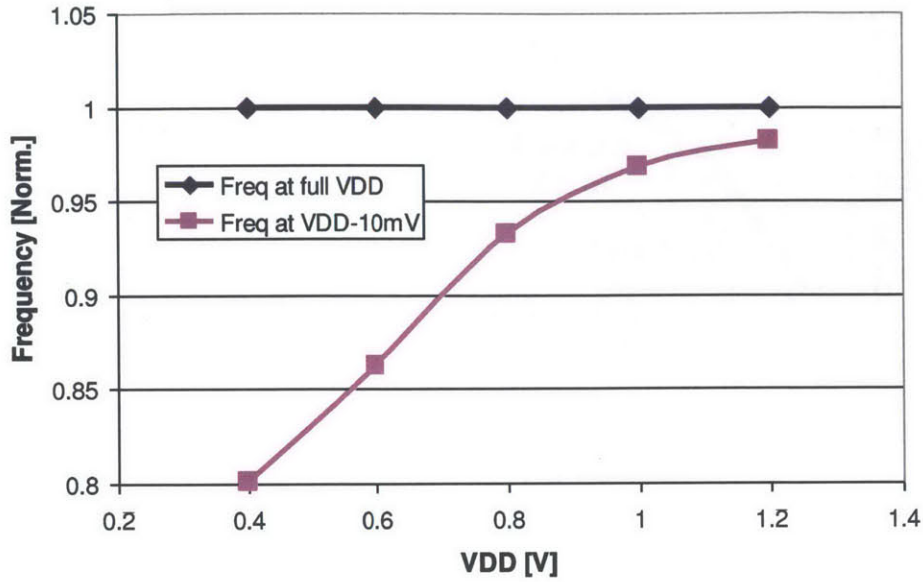


Figure 2-4: The required frequency adjustment for a $\Delta V=10\text{mV}$.

retain their states which cannot be tolerated.

In this system, ΔV can be set with 10 mV step sizes. This value is selected sufficiently large to capture an average EOP across many cycles and small enough in order not to degrade SRAM operation especially at low voltages.

A second consideration is the effect of voltage ripple to the timing. In order not to create any timing related failures, during energy sensing, the circuit under test needs to run with its maximum frequency it gets at $V_{DD} - \Delta V$. This is important to ensure correct operation without any timing errors.

2.2.2 Effect of ΔV Selection On Accuracy

As explained earlier, the energy sensing concept makes the assumption of ΔV being small compared to V_{DD} to make the calculation simpler. However, the actual value of energy lost in the storage capacitor C_{sto} is proportional to $V_{DD}^2 - (V_{DD} - \Delta V)^2$ rather than $V_{DD} \times \Delta V$. Therefore, there is an error introduced in approximating V_{DD} to be equal to $V_{DD} - \Delta V$. A similar analysis was covered in [76]. If it is assumed that $V_{DD} = V_1$ and $V_{DD} - \Delta V = V_2$, the error introduced can be shown as:

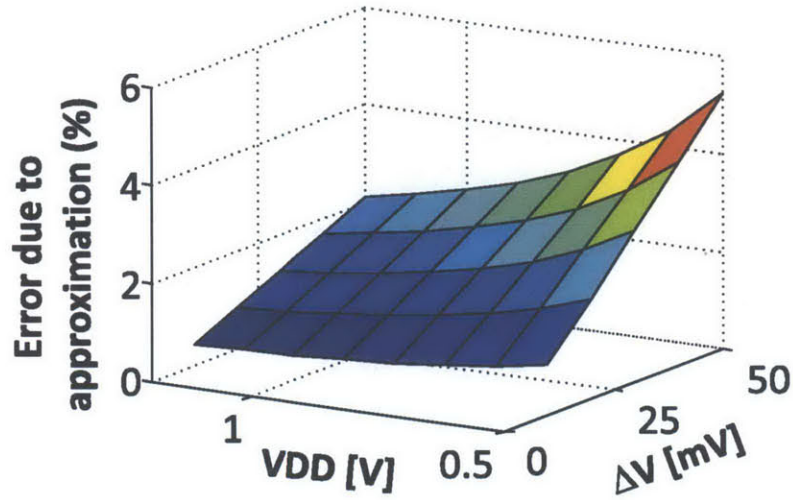


Figure 2-5: Selection of ΔV for accuracy considerations.

$$\frac{\delta EOP}{EOP} = \frac{(V_1 - V_2)}{(V_1 + V_2)}$$

This error is relative to V_1 and V_2 . Especially when V_{DD} is low, ΔV also should be kept low to achieve low error rate in the energy readings. The effect of approximation in energy calculation on accuracy is shown in Figure 2-5. This figure shows how error changes when V_{DD} and ΔV change. At nominal voltage of 1.2 V, even a 50 mV ΔV results into a small ($< 2\%$) error. However at 0.4 V, over 2% error would be introduced if $\Delta V > 20$ mV. Therefore, as observed, Δ needs to be carefully selected to achieve high accuracy in energy readings.

2.2.3 Effect of Comparator Offset on Accuracy

The second major source of error in EOP calculation is the comparator offset. Assuming comparator offset (V_{ofs}), the error introduced in the measure of EOP is given by:

$$\delta EOP_{ofs} = C_{STO}V_1(V_1 - V_2) - CV_1(V_1 - V_2 - V_{ofs}) = C_{STO}V_1\Delta V$$

Therefore, the relative error due to offset is given by:

$$\frac{\delta EOP}{EOP} = \frac{(V_{ofs})}{\Delta V}$$

At a ΔV selection of 10mV, a 1mV comparator offset brings %10 accuracy error to the absolute value of the EOP. However, it is important to note that comparator offset is a static number and it is fixed for every EOP calculation. This means that the absolute error will be the same for each EOP reading. Therefore, for applications where relative error numbers are important, this effect will be more tolerable.

To cope with the comparator error, the comparator is designed using larger transistors. However, sizing gives diminishing returns since $\Delta V_{TH} \propto \frac{1}{\sqrt{WL}}$. Therefore, trimming PMOS transistors are used inside the comparator to achieve an input offset less than 1mV. The comparator design will be explained in the next section.

2.2.4 Effect of Non-Idealities of the Capacitor

Another criterion that needs to be taken into consideration is the selection of the storage capacitor. The effects of non-idealities of this component can result into degradation in the accuracy of energy monitoring.

A capacitor can be modeled as an equivalent circuit of a capacitor with equivalent series inductor (ESL) and equivalent series resistor (ESR). For the storage capacitor values we are interested in, the ESR values change from a few m Ω 's (for multi-layer ceramic capacitor (MLCC)) up to a few Ω 's (aluminum capacitors). Similarly, typical ESL values change from a few nH to tens of nH's. The self-resonance frequency of the RLC circuit happens at:

$$f = \frac{1}{2\sqrt{L \cdot C}}$$

For the frequency range of energy sensing cycles we are interested in (kHz range), parasitic inductor creates a negligible effect. However, the effect of ESR can be significant. The energy sensing cycle using an aluminum capacitor with ESR=8 Ω and with a MLCC capacitor with ESR=80 m Ω are given in Figure 2-6. As it can be seen from the figure, ESR of a few ohms creates an IR drop on the resistor resulting into a sudden drop on V_{DD} . Therefore, for this circuit to be practical, an MLCC capacitor with low parasitics is required. The non-linear junction capacitor which is connected to the C_{STO} is negligible compared to the off-chip capacitor.

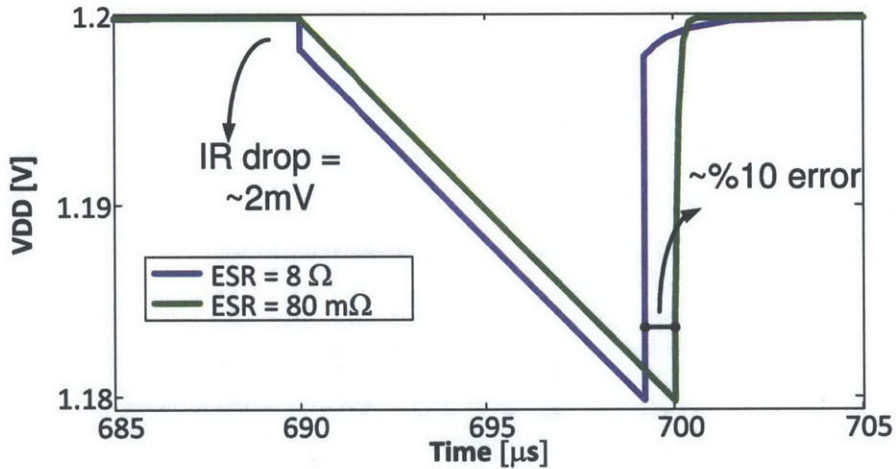


Figure 2-6: Using aluminum capacitor with a large ESR (8 Ω) results into significant error (10%) in energy measurement whereas MLCC parasitics bring negligible error.

Since the measurement of energy assumes a fixed capacitor value at all times, temperature coefficient of the capacitor has a direct effect on the accuracy of the measurement. The temperature coefficient of ceramic capacitors varies drastically depending on the type of the dielectric. The temperature coefficient ($\Delta C/C$) can change anywhere from 30% (class 3) to less than 0.5% (class 1). So, for high accuracy, a class 1 type capacitor is required.

Furthermore, the EOP calculation assumes that the actual value of the capacitor

is known. If the selected capacitor value is comparable to the on-chip and on-board parasitics, the effective capacitor value might need to be measured before the operation. Although not implemented in this work, one possible method is measuring the time it takes for a known current input to charge or discharge the capacitor.

2.3 Energy Sensor Demonstration for an SRAM Application

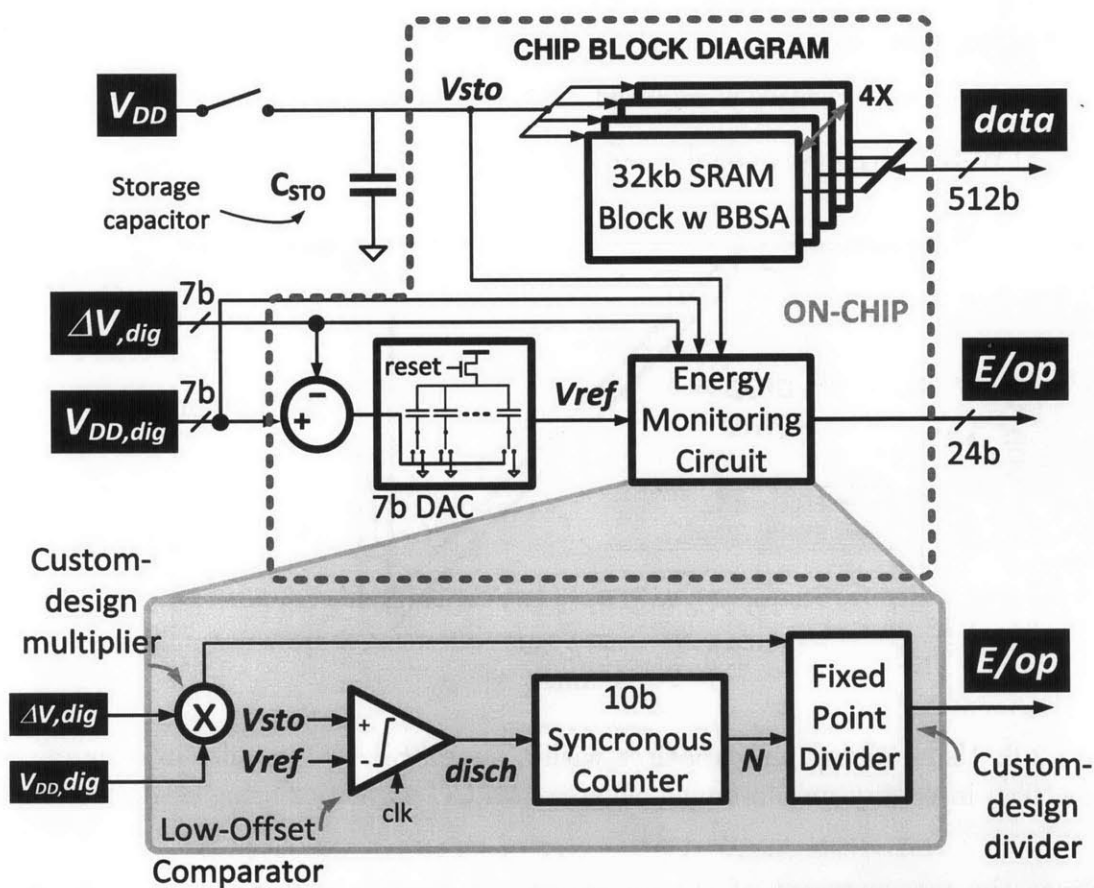


Figure 2-7: Chip block diagram with a focus on the architecture of the energy-sensing circuit.

Figure 2-7 shows the test-chip block diagram with a focus of the energy monitoring circuit. As it can be seen in the figure, the 128 kbit SRAM's power grid is connected to V_{sto} , which is also the voltage across an off-chip storage capacitor, C_{sto} . This node is also connected to an off-chip power supply through a switch. The 7 bit digital

representation of the V_{DD} and ΔV are given as inputs to the system. Although not implemented in this design, the voltage V_{DD} can be generated using $V_{DD,dig}$ with the help of a DC-DC converter in a complete system.

As explained earlier, if during energy sensing, V_{STO} drops from V_1 to V_2 , the EOP can be calculated as:

$$E_{OP} \approx \frac{C_{sto} \times (V_1 \times \Delta V)}{N}$$

The operation of the energy sensing circuit which is shown in Figure 2-7 can be explained follows:

- Before an energy sensing cycle starts, the switch is ON and SRAMs are powered up by the off-chip power supply. Similarly, C_{STO} is charged up to V_{DD} .
- As soon as a new energy sensing cycle is asserted, the analog voltage, V_{ref} , is refreshed using an on-chip 7 bit, capacitive divider type digital to analog converter (DAC). The input to the DAC is the digital word $V_{DD,dig} - \Delta V_{dig}$ which is calculated with the help of an on-chip subtractor.
- Similarly, when energy sensing is activated, the switch turns OFF disconnecting SRAMs from the off-chip power supply, and V_{STO} starts to drop from V_{DD} across C_{STO} . Therefore, during energy sensing cycle, the charge across the C_{sto} powers up the SRAMs. During this time, a 10bit synchronous counter is enabled and it starts to count the number of clock cycles. Meanwhile, an on-chip, low-offset comparator compares V_{STO} to V_{REF} .
- After the N clock cycles, the comparator detects that V_{STO} reaches V_{REF} and stops the counter. At the same time, the switch turns ON and starts to charge the V_{STO} back to V_{DD} .
- Then, EOP is calculated based on N with the help of an on-chip multiplier and divider. When the output is ready, the energy sensing circuit generates a **Done**

signal. Afterwards, the 24 bit digital representation of EOP is outputted as a serial 1-bit output.

- After this phase, the counter is resetted and the energy sensor is ready for the next sensing cycle.

Next, three important blocks (DAC, arithmetic units and the comparator) of the energy monitoring circuit will be explained.

Digital-to-Analog Converter

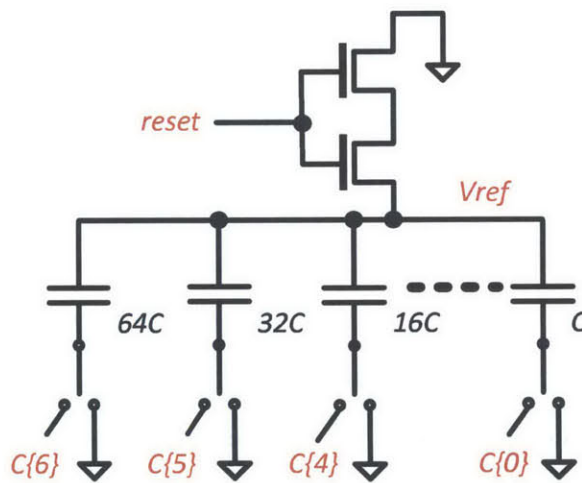


Figure 2-8: Schematics of the DAC.

The reference voltage of $V_{DD} - \Delta V$ is calculated using a 7-bit DAC. A charge redistribution type, capacitive divider DAC is chosen since the DAC needs to consume a very low power. This topology is well suited for such an application. It consumes no static current and very small dynamic current. 7-bit resolution is chosen to enable voltages with 10mV step sizes, which is equal to the ΔV . Even finer implementations can be possible. However, for each added extra bit, the required area almost doubles and the design gets more complicated. Therefore, 7-bit resolution is a good compromise between the required area and resolution.

The schematics of the DAC can be seen in Figure 2-8. The operation of the DAC is as follows. At the beginning of the operation, the top and bottom plates of the

capacitors are shorted to ground to discharge any charge stored on the capacitors. This is performed by deriving C[6:0] to '0' and turning the two NMOS transistors ON by asserting the *reset* input. Then, *reset* is deasserted and V_{REF} node starts to float. Afterwards, the C[6:0] inputs are changed to their required digital value. Thus, the output voltage settles to the analog value of this digital input due to charge sharing between the capacitors.

Due to the leakage, the output voltage of the DAC can change if waited long enough. A cascade of two NMOS transistors is chosen rather than one to decrease the leakage.

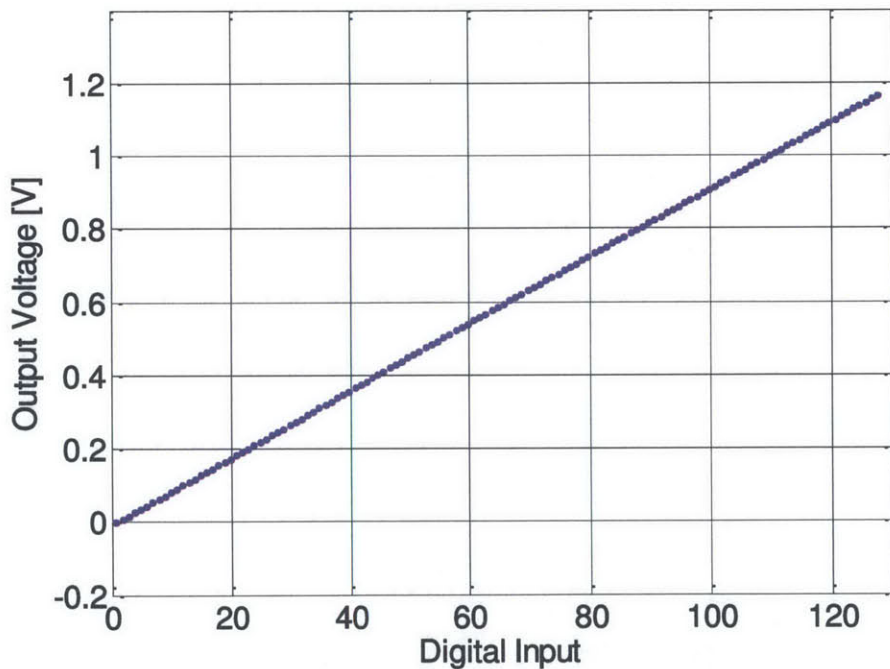


Figure 2-9: DAC output span for 128 different input words.

The analog voltage output, V_{REF} , span that is generated by the DAC can be seen in Figure 2-9. This figure shows the outputs for the full input span from all zeroes to all ones. As it can be seen from the figure, V_{REF} can be generated with less than 10mV step sizes.

Multiplier, Divider and Subtractor Circuits

In this work, the multiplier, divider and subtractor are designed using custom design circuits. They are optimized for minimum area and energy consumption with a trade-off of performance. The designs can be summarized as follows:

1. Subtractor:

In this design, a 7 bit ripple-borrow subtractor with 1 bit full subtractors is used. This design only requires implementing a full subtractor block and replicating it 7 times. An example 4 bit ripple-borrow subtractor can be seen in Figure 2-10. The advantage of this subtractor is that its design is simple and it requires a small area and energy consumption.

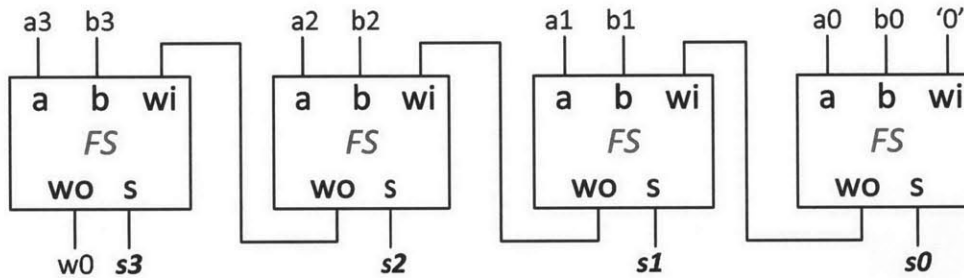


Figure 2-10: A 4 bit ripple-borrow subtractor using full subtractors.

2. Multiplier:

In this work, the multiplier is needed for calculating $C_{STO} \times V_{DD,dig}$ and ΔV . Both are represented as 7 bit digital numbers and therefore, the output of the multiplier is a 14 bit digital number. When designing multipliers, there is always a compromise to be made between how fast the multiplication process is done versus how much hardware is required. Since, in this energy sensor, the area and energy consumption are important, the shift-and-add method is chosen. This is a simple multiplication method that is slow but efficient in use of hardware. This way, the multiplication takes 7 clock cycles. This is tolerable for the energy sensor implementation since the sensing operation is likely to be performed to re-optimize the system due to changing system dynamics which

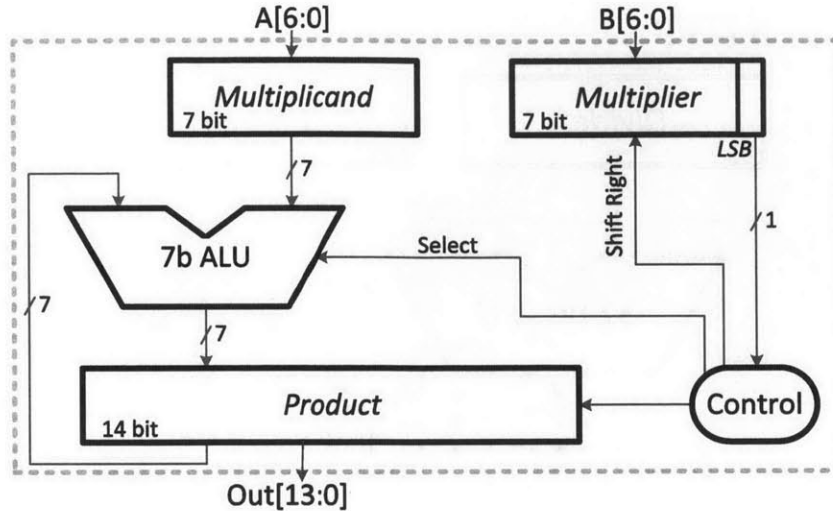


Figure 2-11: Block diagram of the custom design shift-and-add multiplier.

will not be happening within hundreds of clock cycles. Although this topology might bring a larger energy overhead, it results into a small area. In this design, we chose to minimize the area consumption for the multiplier and the divider circuits.

The block diagram of the multiplier can be seen in Figure 2-11. Shift-and-add multiplication is similar to the multiplication performed by paper and pencil.

3. Divider:

The 14 bit digital representation of $C_{STO} \times V_{DD,dig} \times \Delta V$ needs to be divided by the 10 bit digital number N which represents the number of clock cycles. Therefore, the result is a 24 bit number. Similar to the multiplier, it is designed for minimum area and power consumption. Therefore, a shift-and-subtract divider is chosen. This is a simple divider which trades off performance for better area and energy. Due to that reason, the division takes 14 cycles. Due to the same reasoning explained in multiplier section, this is a tolerable trade-off.

The block diagram of the divider can be seen in Figure 2-12. In this implementation, at every clock cycle, the divider is shifted to the right. It is subtracted from the dividend. If the result is negative, '0' is shifted to the quotient. Other-

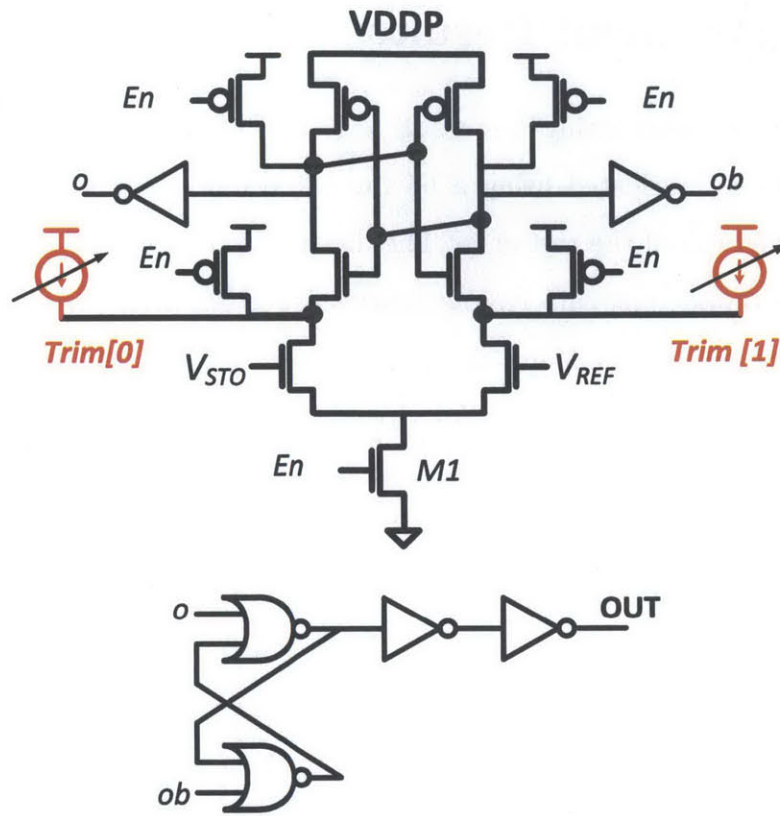


Figure 2-13: Trimmed strong-arm type sense amplifier is used for comparator.

The schematics of the comparator can be seen in Figure 2-13. An offset below 1mV is achieved thanks to the current trimming and transistor sizing.

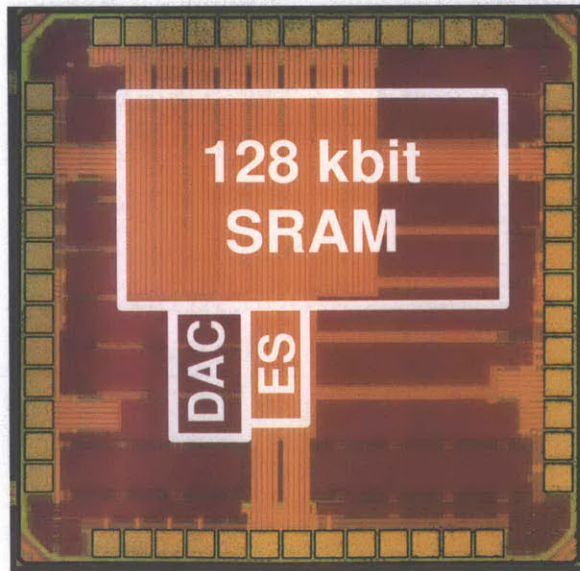


Figure 2-14: Die photo of the 128 kb SRAM in 65 nm CMOS.

2.4 Measurement Results

The idea of energy monitoring for a 128 kbit SRAM is illustrated in a test-chip prototype which is fabricated using a 65 nm LP CMOS process [72]. Figure 2-14 shows the micrograph of the test-chip. The die size is 1.4 mm by 1.4 mm.

Some of the important chip characteristics are given in Table 2.1. SRAMs are measured to achieve functionality down to $V_{DD} = 0.37$ V thanks to 8T bit-cells and write assists. A pseudo-differential strong-arm type SA is used for sensing and one of the inputs is connected to a 100mV reference voltage. The test-chip also proposes a sense amplifier offset compensation technique using body biasing of the input transistors with a voltage of $V_{DD} + 150$ mV. The design of the SA and the SRAM will be explained in more detail in Chapter 4.

The area of the ES is 0.026 mm^2 and the DAC is 0.059 mm^2 . Together, they bring a 16% area overhead compared to the SRAM area. However, for a system application with larger SRAM sizes, this area overhead will be smaller.

Table 2.1: Test-chip specifications for the 128kb SRAM with an embedded energy monitoring circuit.

Technology	65 nm LP CMOS
Chip Size	1.4 mm x 1.4 mm
Configuration	4 blocks of 256 rows, 128 columns
DAC resolution	7 bits (10 mV steps)
Supply Voltages	DAC: 1.2 V VDD: 0.37 to 1.2 V VDDH: 0.6 to 1.2 V VDDDB: AVDD + 150 mV REF: 100 mV
Area	DAC: 0.059 mm^2 ES: 0.026 mm^2 SRAM: 0.524 mm^2
EOP	108 pJ at 1.2 V 49.2 pJ at 0.8 V 21.2 pJ at 0.5 V

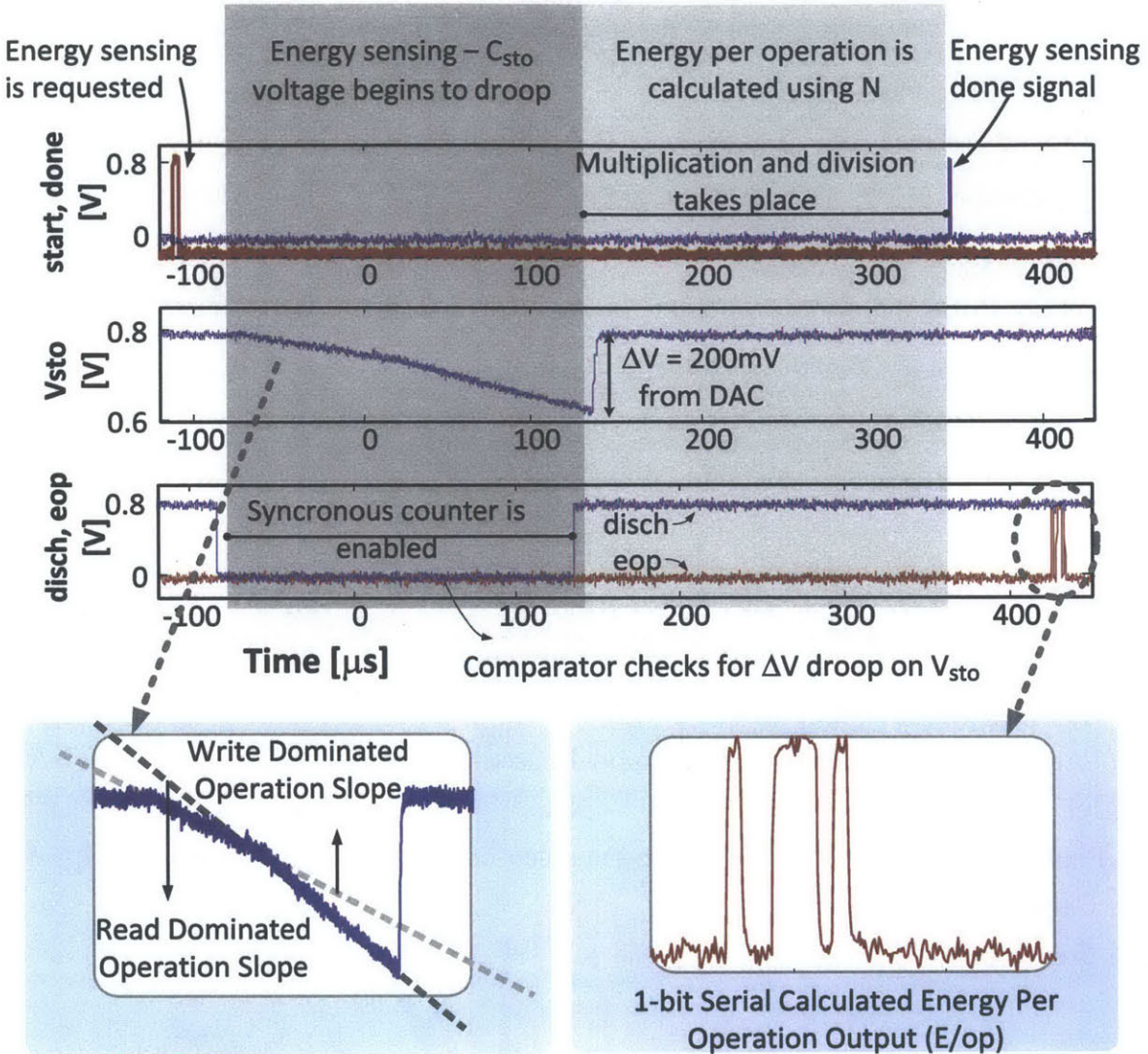


Figure 2-15: Oscilloscope outputs for the critical signals of energy sensing circuit.

Figure 2-15 demonstrates the operation of the energy sensing circuit by showing the oscilloscope outputs for its critical signals. An energy sensing cycle is initiated by asserting ESstart. Then, a custom design 7-bit capacitive DAC generates the reference voltage, $V_{ref} = V_{DD} - \Delta V$. The comparator output (disch) keeps a 10-bit synchronous counter enabled during this period to generate N. When V_{sto} drops below V_{REF} , a custom-design multiplier and fixed-point divider are used to calculate absolute EOP. The EOP is outputted as a single bit serial output at the end of energy sensing operation.

In this design, the write operation is more energy consuming compared to read

which means that during a write dominated operation, supply voltage value drops with a steeper slope.

Table 2.2 shows the comparison of the proposed energy sensing circuit to previously published work. This implementation achieves 10% accuracy without a calibration process requirement. The current implementation requires an off-chip capacitor for energy storage, however, the filtering capacitor of a power converter can be used as C_{sto} to decrease the number of off-chip components required. The proposed monitor does not let the voltage drop arbitrarily during sensing which is important especially for an SRAM application for maintaining data retention.

Table 2.2: Comparison of the proposed energy sensor with recent work.

	This work	[22]	[33]	[76]
Technology	65 nm CMOS	32 nm CMOS	NA	65 nm CMOS
V_{DD} [V]	0.34 - 1.2	0.7 - 1.15	3.3	0.25 - 0.7
Accuracy	10%	model based	20%	NA
Area Overhead	16%	NA	NA	21%
Requirement	capacitor	thermal sensor	PFM mode DC-DC	capacitor
Calibration Required?	No	NA	Yes	No

2.4.1 Measured EOP Accuracy

Figure 2-16 shows the measured EOP numbers on the y-axis and the V_{DD} on the x-axis. EOP values are scaled down with voltage scaling by more than $5\times$ which shows the benefit of voltage scaling on energy efficiency.

Another important point of this figure is the fact that it shows the accuracy of the EOP calculation. This figure shows the EOP calculation using three different methods.

1. Reading the digital EOP output of the energy monitoring circuit.
2. Measuring energy consumption values of the test-chip.

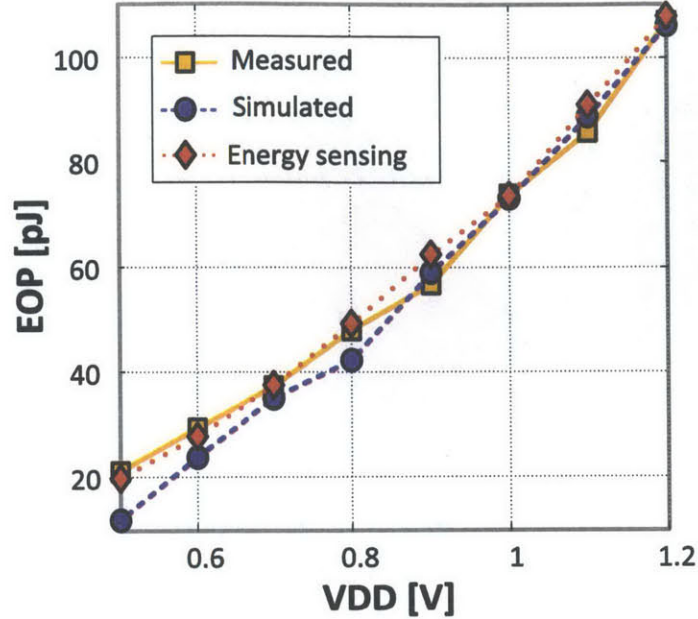


Figure 2-16: EOP vs. V_{DD} graph using three methods: 1- measured, 2- simulated, and 3- energy monitoring circuit output.

3. Using transistor-level extracted simulation with parasitics.

As it can be seen from the figure, for the V_{DD} range of 1.2 down to 0.5 V, EOP generated by the energy sensor is within 10% range of the measured value. This shows that the energy sensor can generate accurate results.

2.4.2 EOP Under Dynamic Effects

It is important to observe how EOP changes with dynamic effects. In this context, Figure 2-17 shows measured EOP calculated by the on-chip energy sensing circuit across different temperatures (from 30°C to 80°C) and by varying total read and write operation ratios (0 to 1). Read accesses result in smaller energy consumption and the ratio of total read and write operations can significantly change from one application to another. Moreover, higher temperature results in larger leakage and larger overall EOP. Depending on those conditions, EOP can change by more than $2\times$ as shown in the figure.

This example illustrates only a portion of the dynamics a complex system can have. Under the effect of multiple dynamics (such as application load changes, tem-

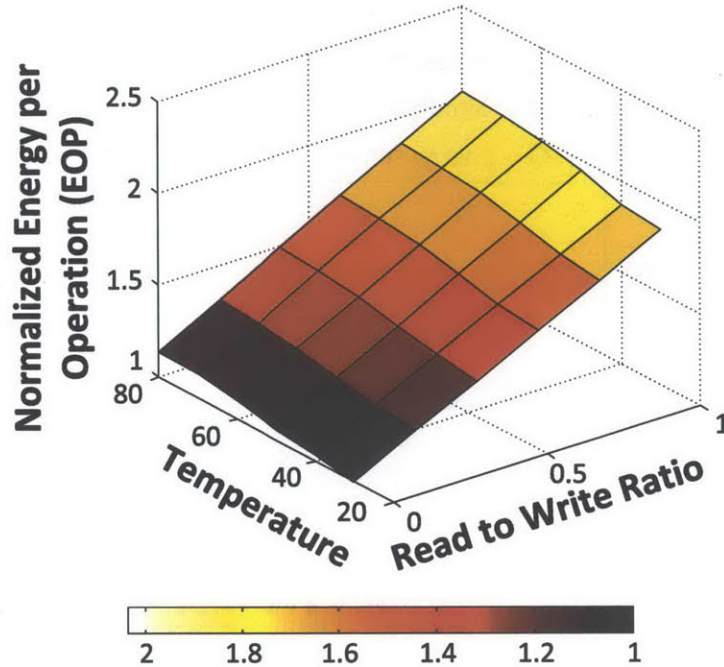


Figure 2-17: Measured EOP under different read operation to write operation ratios and temperatures.

perature fluctuations, process variations, aging effects, and voltage fluctuations), it is hard to model SRAM energy consumption at the design phase and those effects can significantly change the system optimization. Therefore, systems can greatly benefit from run time energy sensing calculations of SRAMs.

2.5 Summary and Conclusions

Today's complex computing systems need to continue working efficiently under dynamic operating conditions and competing design targets. To enable on-fly power and performance optimizations, recent systems leverage power management engines that use monitoring circuits based on energy models. However, models cannot fully represent the actual profile of a complex system. In this context, energy sensing circuits that can provide actual energy consumption information would be useful for system level optimizations.

SRAMs are one of the fundamental building blocks of today's complex systems since they account for a large portion of total energy consumption, area and access

time. Therefore, adaptive systems with embedded SRAMs can greatly benefit from run-time SRAM energy consumption information for better system level power and performance optimizations. This chapter proposed using an energy sensing circuit that is capable of generating a digital representation of the absolute EOP consumed by SRAMs.

In this chapter, the challenges and design considerations of the energy monitoring circuit are investigated. Firstly, the monitoring circuit needs to be energy-efficient so that the system can get net energy savings by using it. Secondly, it needs to be non-intrusive to the operation of the circuit under test. This way, it can be designed independently. Also, its area needs to be small, especially for applications that can benefit from using multiple energy monitors. Lastly, it needs to provide accurate results. The error sources that affect the accuracy of the sensor are analyzed.

A prototype test-chip is designed using a 65nm CMOS process to demonstrate the energy sensing operation for a 128kb SRAM application. In this chapter, the building blocks of the monitor are investigated. Some of those blocks are the comparator, arithmetic structures and DAC.

The measurement results of this test-chip show that the sensor can generate results within 10% of the actual energy consumption. Furthermore, depending on changing dynamic conditions (temperature and application type), the EOP of the SRAM can change more than $2\times$. The energy monitor required a 16% area overhead but this number is expected to be smaller for larger SRAM sizes. The dynamic power overhead was $<1\%$. These specifications of the monitoring circuit motivate its usage for system-level power management.

Chapter 3

Self-Aware Processor Design Using Embedded Energy Monitors

Angstrom is an MIT-led project [77] with the goal of creating technologies necessary for extremely scaled computers. In this context, Angstrom targets massively many-core processor systems (up to 1000-cores). Those massively scaled computers face several major challenges such as energy efficiency and programmability. Angstrom processor addresses those challenges by supporting self-aware computing models, systems and circuits. Figure 3-1 shows the block diagram of the Angstrom processor, and as it can be seen, it is a fully distributed, scalable design with voltage scalable and adaptive hardware structures.

For the Angstrom processor, a self-aware computation model, called *SEEC*, is proposed as the software decision mechanism. SEEC has been conceptually shown to alter the behavior of a system to meet multiple goals and automatically adapt to environmental changes [32]. Simulations of a 256-core Angstrom system show that exposing hardware adaptation to the software management system has the potential to improve performance per watt by an average of over 100% compared to a non-adaptive system [20].

Like all self-aware systems, SEEC is characterized by the presence of an *observe-decide-act* (ODA) loop [28, 29]. It needs to continuously monitor its goals (*observe*), and it needs to be aware of its available resources (*actions*). This way it can determine

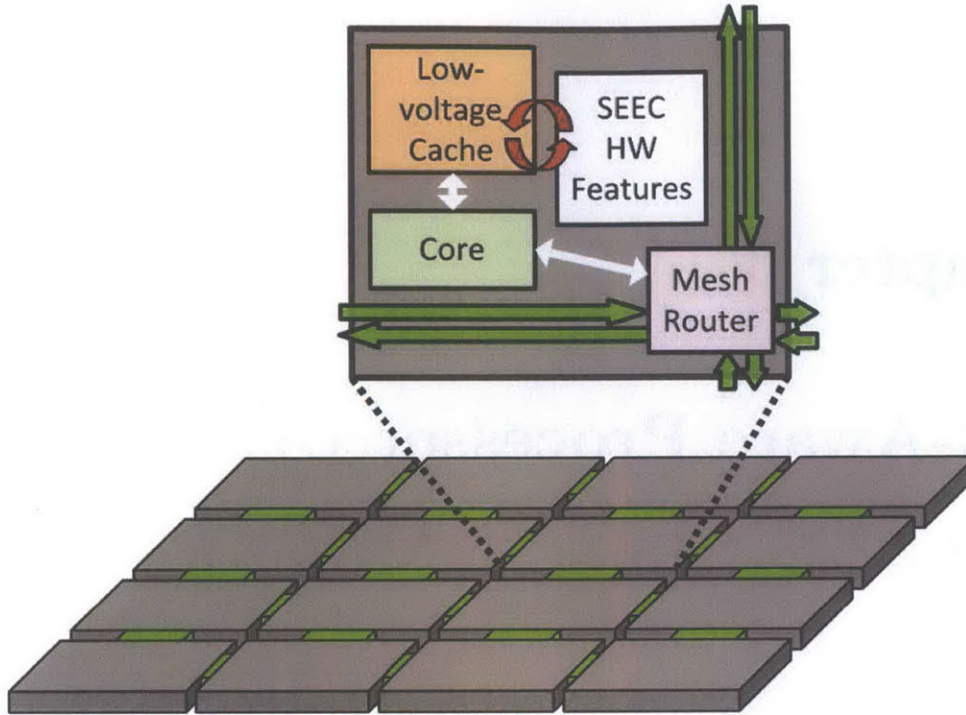


Figure 3-1: Angstrom multicore architecture.

how best to use its resources to meet goals (*decide*). Obviously, to enable such a system in the hardware, systems need (1) building blocks that can enable observation of important system metrics such as energy (or power) consumption and performance (2) reconfigurable system knobs to adapt the system to changing dynamics.

In recent complex processor systems, a popular way to monitor power consumption is to use blocks that can estimate it based on power models. However, models cannot fully represent the actual profile of a complex and dynamic processor system [22]. An absolute energy monitoring circuit is demonstrated in [72], as explained in Chapter 2, but additional benefits can be obtained by integrating them within the DC-DC converters.

This chapter presents a self-aware processor system, which is designed as a single-core example of the Angstrom processor. This design uses multiple energy monitors that measure the actual energy consumption of important blocks of the system. Furthermore, it utilizes reconfigurable and voltage scalable circuits for adaptation. The system is designed to work with SEEC to decide how to best allocate the resources

based on both hardware and software conditions and constraints. A prototype system is designed using a 0.18 μm CMOS technology.

3.1 Observe-Decide-Act Loop For Self-Aware Processor

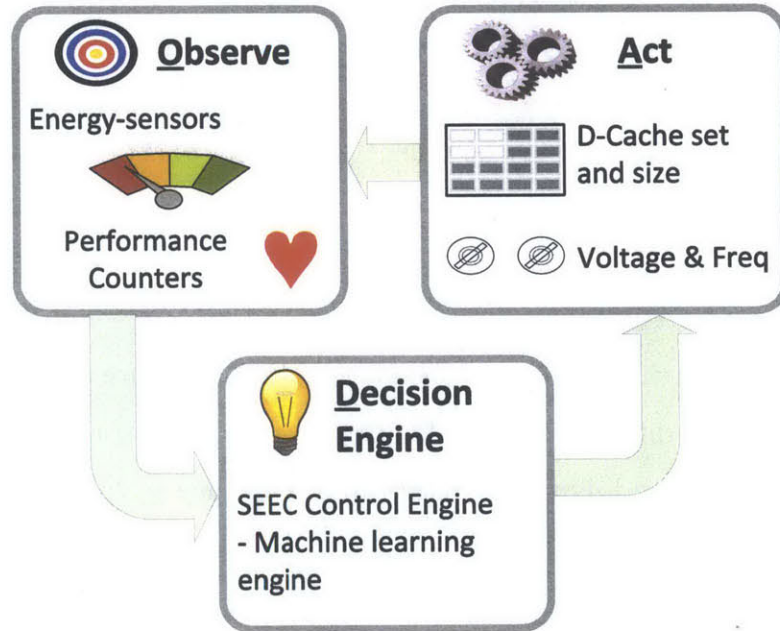


Figure 3-2: ODA loop for the self-aware processor.

The self-aware processor system proposed in this thesis supports *open computing* by exposing observations and action knobs to the software runtime decision engine. The self-awareness of the processor can be characterized by an ODA loop where different components of the system contribute to *1-observation*, *2-action* and *3-decision*. This ODA loop is illustrated in Figure 3-2 and its specifics can be explained as follows:

1. Observation

In the self-aware processor system, the observations are performed by both traditional performance counters and the proposed energy monitoring circuits. The performance counters are useful since they provide valuable insight about the behavior of an application on a particular hardware. In this self-aware

processor system, they can count memory operations, cache hits and misses, and pipeline stall cycles.

In addition, the self-aware processor system has multiple energy monitoring circuits. Those circuits adopt the same energy sensing concept explained in Chapter 2. However, the ones in the self-aware system are embedded into on-chip DC-DC converters. This way, they share the filtering components of the DC-DC converters and do not require extra off-chip components. (The filtering capacitor is still off-chip). Similarly, they use the multiplier and divider in the core pipeline which makes their design simpler.

2. Action

In order to alter the behavior of the system, the self-aware processor is exposed to a number of different actions knobs. In this system, those are data-cache reconfigurability and dynamic voltage and frequency scaling.

3. Decision

The self-aware processor system is designed to use SEEC as its runtime decision engine. SEEC can alter the system components for better power and performance optimization. The SEEC runtime system will often have to make decisions about actions and applications with which it has no prior experience. In addition, the runtime system will need to react quickly to changes in application load and fluctuations in the available resources. To meet those requirements for handling general and volatile environments, SEEC engine is designed with multiple levels of adaptation. At the lowest level, SEEC acts as a classical control system, taking feedback in the form of heartbeats and using it to tune actuators to meet its goals. Additional layers of adaptation, including hardware control and machine learning based techniques allow SEEC runtime to allocate resources efficiently without prior knowledge of the application [20].

3.2 LEON3 Processor

The self-aware processor core is designed based on a LEON3 core and the important features of this core are given in Appendix C. This section will give a brief description about the modifications performed on the original LEON3 architecture (vanilla LEON3).

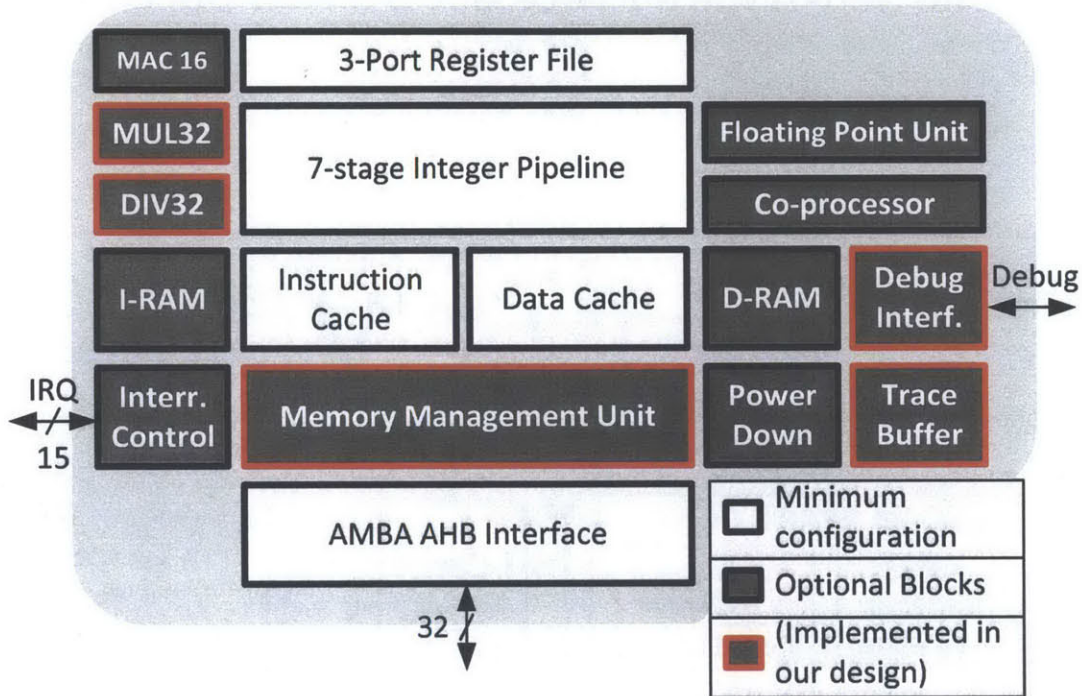


Figure 3-3: LEON3 processor core block diagram.

The block diagram of the vanilla LEON3 can be seen in Figure 3-3. It comes with some mandatory blocks and some optional blocks. The blocks that are shown in red are the optional blocks that were chosen to be implemented in the self-aware system.

In the self-aware processor, in order to achieve a high-density solution, we replaced the register-based caches with custom-design SRAM circuits. This resulted into more than $4\times$ cache area reduction. The area savings can even be larger for designs using SRAMs with industry sized bit-cells.

In the self-aware processor, the multiplier and divider blocks of the vanilla LEON3 are adopted. On the other hand, our processor do not implement an FPU block since if implemented, it would have resulted into more than doubling the total gate count.

Instead, self-aware processor support integer based operations. Our processor also does not support MAC instructions, and does not have a co-processor. Furthermore, it does not adopt the power-down mode since for the single-core self-aware processor, separate blocks can be powered down rather than the entire core.

3.3 The Self-Aware Processor System-on-Chip Implementation

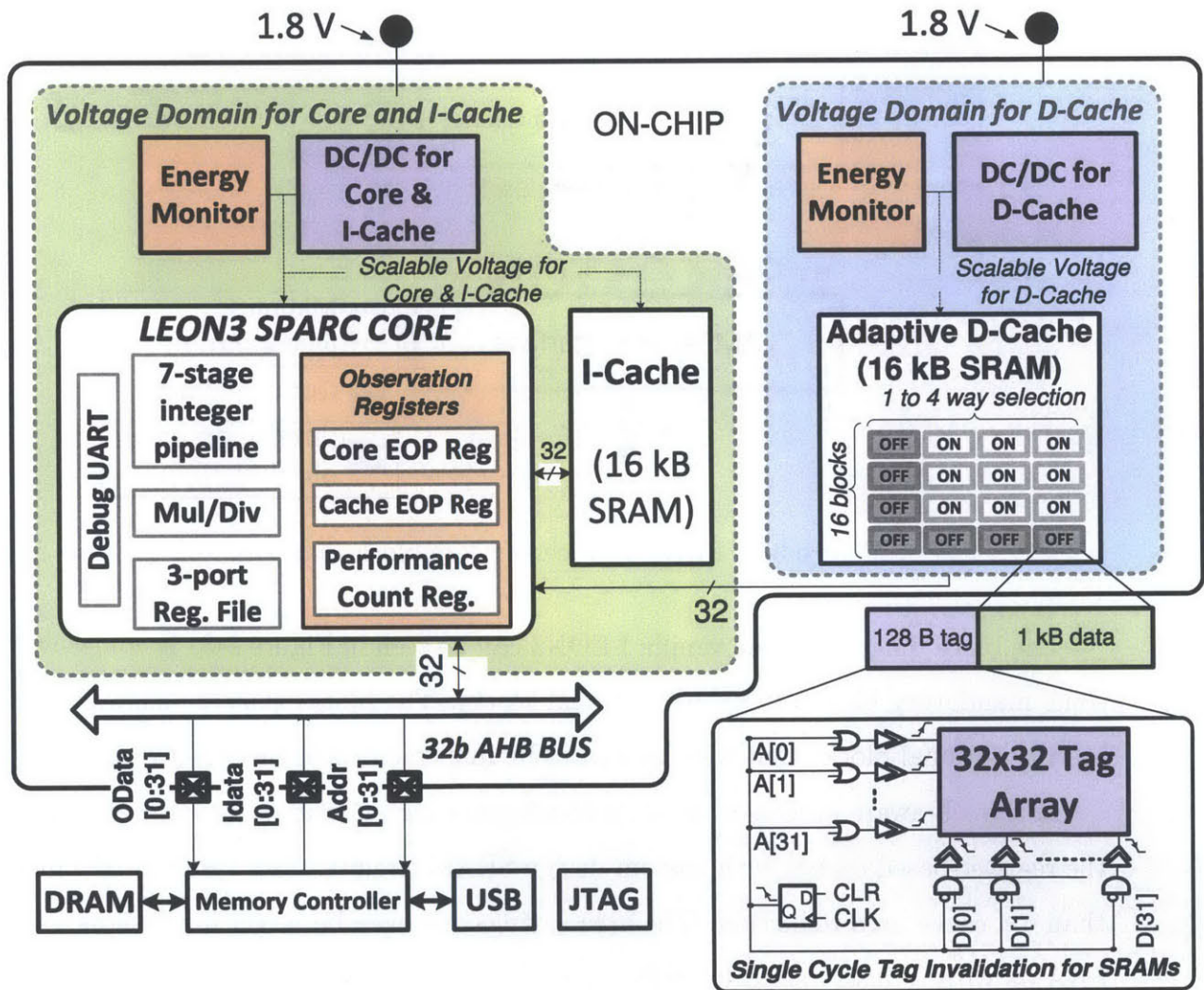


Figure 3-4: The self-aware processor block diagram.

Figure 3-4 shows the block diagram of the self-aware processor system. This

system takes advantage of the highly reconfigurable structure of the LEON3 model and a modified LEON3 single-core processor constitutes the core of the system.

The self-aware processor system would benefit from the capability of distinguishing between the energy spent on computational operations versus energy spent on data storage. Thus, in this processor system, the core and the instruction-cache (*i-cache*) operate from one voltage domain; whereas, the data-cache (*d-cache*) operates from another voltage domain. Those two voltages can be separately controlled and scaled. This way, the system can optimize its configuration better for cache-bound vs. compute-bound applications independently.

Operating digital circuits under voltage scaling necessitates highly efficient DC-DC converters that are capable of delivering a wide range of voltages from the nominal supply voltage. In the self-aware processor system, two DC-DC converters are designed to deliver variable load voltages from 0.6V to 1.8V using a 1.8V supply. Those two converters power up the d-cache and the core domains separately.

Those two domains not only can be voltage scaled independently, but also their energy can be monitored separately. The principle of energy monitoring that is used in the energy sensors is similar to the concept that was proposed in Chapter 2. But the two separate energy monitors are embedded into the two DC-DC converters. During energy monitoring, the number of clock cycles it takes for a ΔV voltage to appear on the supply voltage (which is also known as number N) is stored in the *Core EOP Register* and *Cache EOP Register* respectively. These registers can be seen in Figure 3-4. Then, using the pipeline multiplier and divider, EOP can be calculated using the digital representation of the voltage, capacitor, ΔV , and N .

In this system, i-caches, d-caches and the trace buffer are designed using custom-built SRAM circuits. However, traditional SRAM circuits with 6T bit-cells cannot enable voltage scaling from 1.8V to 0.6V due to the degraded noise margins at low voltages. Therefore, assist techniques and different bit-cell topologies might be required. In this design, 8T bit-cell based SRAMs with write-assists are used in order to enable an operation down to 0.6V. The design and analysis of a similar 8T bit-cell based, low-voltage SRAM (which is designed in 65nm) will be explained in Chapter

4.

For a self-aware system, having the d-cache to be reconfigurable is very beneficial for reducing the power consumption for the same performance. Disabling unnecessary parts of the d-caches helps the processor to optimize power and performance trade-offs. In this context, associativity and size reconfigurability enable different trade-offs and have been used in reconfigurable processor systems extensively [22, 23]. For instance, for a small working set that can fit into a smaller memory size, it would help to reduce the cache size to achieve overall power reduction. However, for another application, increasing the cache size or associativity can be more desirable. To enable this, in this design the custom SRAM memories are designed to be set and way associative.

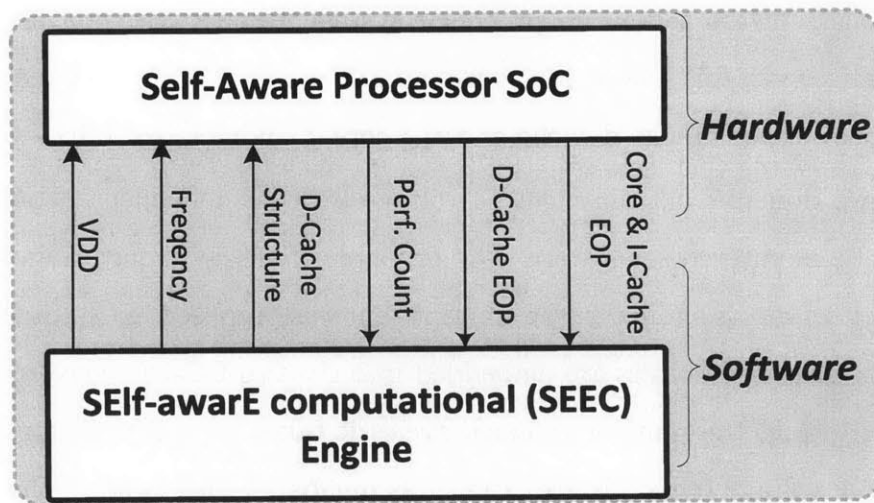


Figure 3-5: Hardware and software separation in self-aware processor system.

The hardware and software separation of the self-aware processor is illustrated in Figure 3-5. The SEEC decision engine resides in the software as library specific functions. On the other hand, the energy sensors and performance counters, which are used for observation, reside in the hardware. The configurable circuits like the d-cache and voltage and frequency scaling capabilities also reside in the hardware. However, being an *open* system, the self-aware processor makes all those hardware components accessible by the software. Similarly, the hardware reconfigurability can be controlled by the software.

3.4 DC-DC converters with Embedded Energy Monitors

Operating digital circuits at a wide voltage range including near threshold region necessitates a low power and high efficient DC-DC converter. The converters in the self-aware processor need to be able to deliver the wide voltage range of 1.8V down to 0.6V from a 1.8V supply.

There are various DC-DC converter topologies with different trade-offs. The most popular three are the linear regulators, switched capacitor converters and switching regulators. The modulation mechanisms and efficiency calculations of switching regulators are briefly summarized in Appendix D.

This chapter focuses on the DC-DC converter implementation used in the self-aware processor which are designed using PFM mode buck-converters.

3.4.1 DC-DC Converter Implementation for the Self-Aware Processor

In the self-aware processor, the variable supply voltages for the two domains are created using two separate DC-DC converters. The load power of one of the domains can be seen in Figure 3-6. The load power varies between 150mW to 0.8mW as the voltage scales. While a PWM DC-DC can be made efficient at full load, as the load scales down, PWM control is expected to become less efficient. The advantage with the PFM mode is that the power transistors are switched ON only when necessary (when the load voltage falls below the reference). Therefore, for this application, a PFM mode converter is chosen. Furthermore, the PFM mode modulation is more suitable to work with the energy sensor since the inductor current goes to zero every switching cycle and energy sensing can be started at this period. Thirdly, the control circuitry of the PFM mode converter is much simpler compared to the PWM mode DC-DC converter. By using a PFM mode DC-DC converter, all the necessary blocks, except the filtering components, are designed to be on-chip.

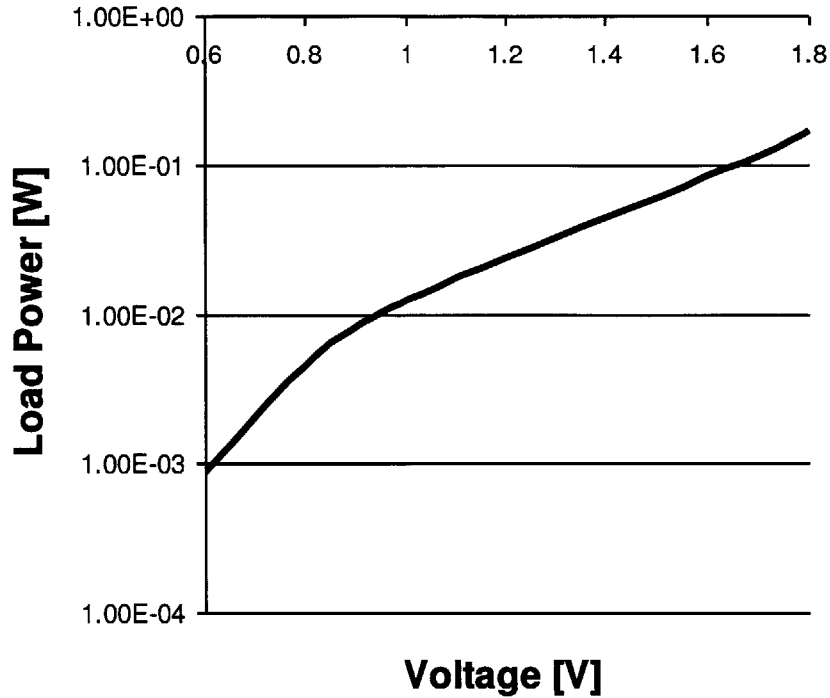


Figure 3-6: The variation of the load power vs. V_{DD} .

Figure 3-7 shows the block diagram of the DC-DC converters with the embedded energy monitoring circuits.

The key components of the design are the filtering transistors (M1 and M2), filtering components (inductor and capacitor), pulse generators, and the reference voltage generator (5b DAC). The designs of these components are explained briefly below.

Power Transistors

As explained in Appendix D, the conduction losses are inversely proportional to the width of the power transistors whereas the switching losses and leakage losses are directly proportional to the width. In $0.18\mu\text{m}$ technology, where the converter is built, the leakage component is negligible since its effect is less than 1%. Therefore, the power transistors trade-off switching loss to the conduction loss; and in this design, power transistor widths are chosen such that the sum of the two losses is minimized.

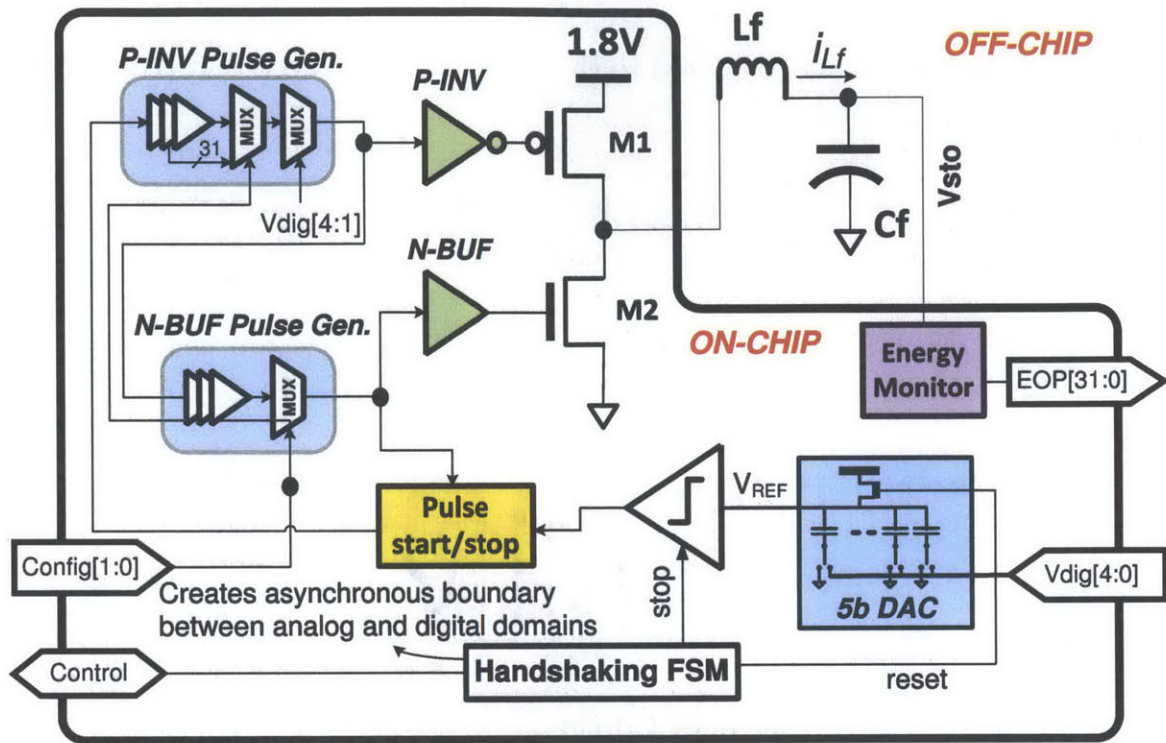


Figure 3-7: The block diagram of the DC-DC converters with the embedded energy monitoring circuit.

Filtering Components

The filtering elements are chosen based on the allowable output ripple and the maximum load current. Additionally, in the DC-DC converter design, the voltage ripple needs to be much smaller compared to the ΔV step size of the energy sensors (50mV) since the ripple directly affects the accuracy of the energy sensing operation. Therefore, the voltage ripple is limited to less than 1mV. The equations used for the filtering capacitor selection are summarized as follows [75]:

$$L_f = \frac{T_{PMOS} \times (V_{in} - V_{out})}{2 \times I_{max}} \quad (3.1)$$

$$C_f = \frac{Q_t}{V_{ripple}} \quad (3.2)$$

where Q_t is the total charge delivered to the capacitor in one cycle and given by the formula:

$$Q_t = \frac{T_{PMOS}^2 (V_{in} - V_{out}) V_{in}}{2V_{out} L_f} \quad (3.3)$$

Pulse Generation Circuit

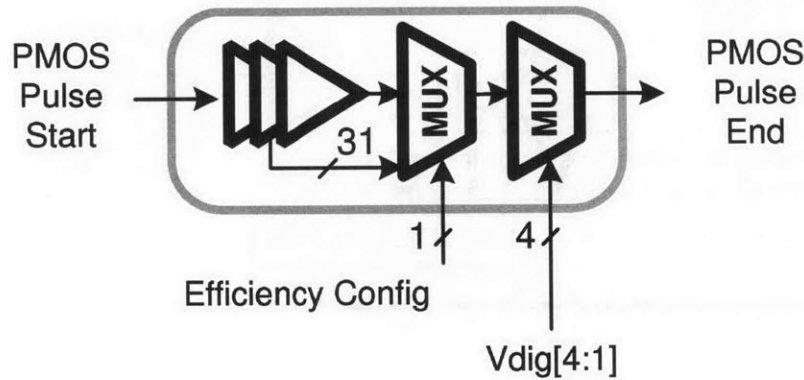


Figure 3-8: PMOS pulse (T_{PMOS}) generation circuit.

The pulse generation circuit for both the NMOS and the PMOS are implemented with the help of delay lines. The delay is obtained using cascaded inverter cells as shown in Figure 3-8. To increase the delay, the inverter lengths are chosen to be larger than the minimum. In the design, the width of the NMOS transistor pulse is fixed at a width, T_{NMOS} . The advantage of fixing T_{NMOS} rather than T_{PMOS} is that the inductor peak current is decreased for lower V_{OUT} values where the delivered output power will also be lower due to low-voltage operation. This way, efficiency is kept above 90%. $V_{dig}[4 : 1]$ signal selects the pulse width required at the operating voltage of 1.8 V to 0.6 V. On the other hand, one bit *Config* signal is used to select a smaller delay for both the PMOS and the NMOS. Changing those two pulse widths together by the same ratio results into an efficiency curve with a different slope as will be shown in Figure 3-14.

The ratio of required NMOS pulse width to PMOS pulse width with the change in V_{OUT} is shown in Figure 3-9. Fixing T_{NMOS} , the width of the PMOS needs to be:

$$T_{PMOS} = \frac{V_{out}}{V_{in} - V_{out}} \quad (3.4)$$

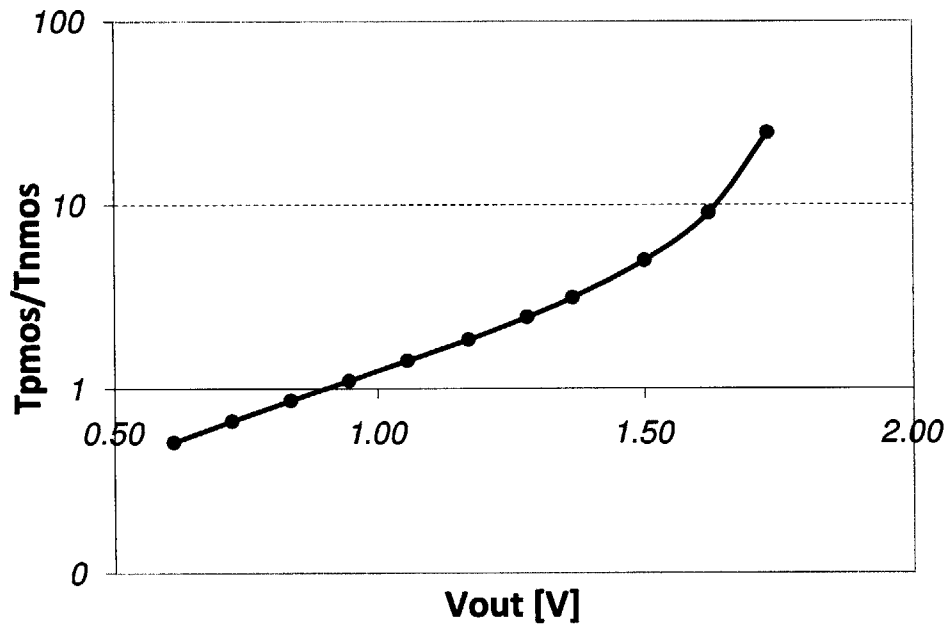


Figure 3-9: Ratio of required PMOS pulse width to NMOS pulse width vs. operating voltage.

Reference Voltage Generator

The operation of the PFM mode DC-DC converter is based on keeping the output voltage higher than a reference voltage. This reference voltage, V_{REF} is an input to a comparator that stops the pulse generation unless the voltage drops below the V_{REF} . This reference voltage generator is performed using a 5 bit DAC. A charge redistribution type, capacitive divider DAC similar to the one explained in Section 2.3 is used.

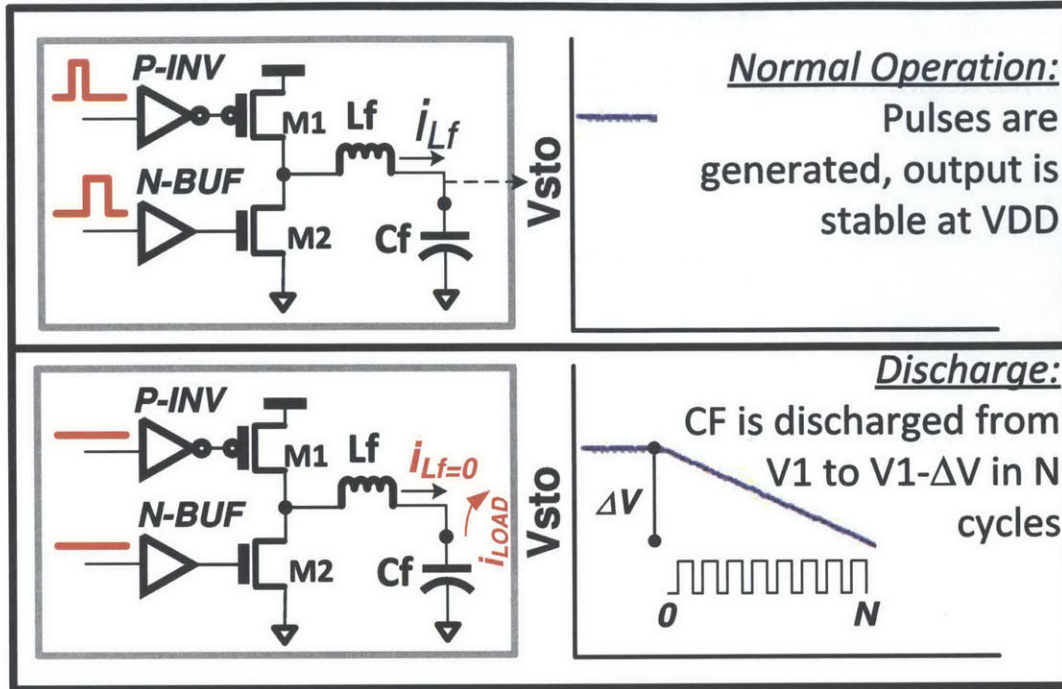


Figure 3-10: The operation of the energy monitors and DC-DC converters together.

3.4.2 Energy Monitoring Circuit Operation with DC-DC

The operation of the energy monitoring circuit is similar to the one shown in Chapter 2. However, in this chapter it is integrated into DC-DC converters and their operations need to be in harmony. The operation of the DC-DC converter with the embedded energy sensor is illustrated in Figure 3-10 and it can be explained in three steps:

1. During *normal operation*, depending on the desired output voltage, the necessary pulses are supplied to M1 and M2.
2. When an energy monitoring period occurs, the monitoring circuit uses a two-step process to generate EOP. Step 1 is the *discharge* phase where the M1 and M2 pulses are kept OFF. The number of clock cycles it takes for a ΔV voltage drop across a known filtering capacitor, C_f , is observed.
3. Then, in step 2, the voltage is restored and EOP is calculated by using $EOP = C_F \times VDD \times \Delta V/N$ assuming ΔV is small.

In this design, ΔV can be set with around 50 mV step sizes through an on-chip capacitive DAC to achieve high monitoring accuracy. On the processor side, dedicated registers are used to adjust the operating voltage of the two domains and to issue energy monitoring operations.

Handshaking Between Core Domain and Fixed Domain

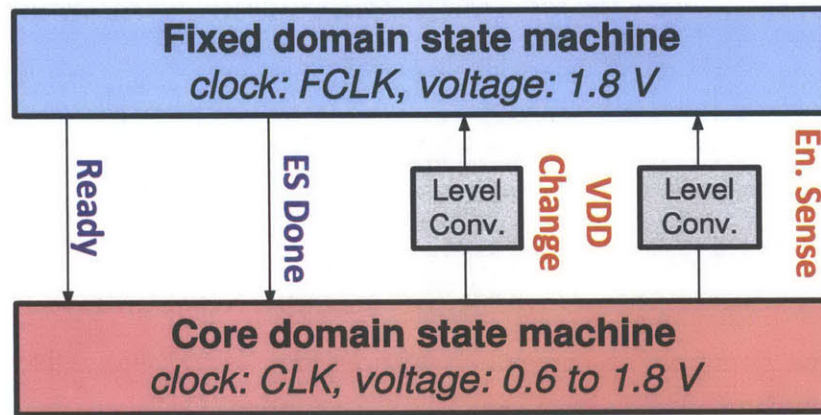


Figure 3-11: The asynchronous boundary block diagram.

One of the challenges of this design was that the DC-DC converter control circuits and energy monitors are designed to work with a fixed clock, whereas the core frequency is adjusted depending on the operating voltage. The fixed clock is designed to be 30 MHz whereas the variable CLK is designed to vary between 1MHz and 100MHz which is handled by the core. Hence, between the core domain and energy monitors, a handshake protocol is created. The block diagram of this domain can be seen in Figure 3-11. Since the core domain voltage can vary between 1.8 to 0.6 V, its output signals, *Energy_Sense* and *VDD_Change*, require to be level shifted to 1.8V. Two DCVSL type level converters are inserted between these two domains.

The important signals of the asynchronous boundary are shown in Figure 3-12. The operation of the handshaking for a voltage change operation is summarized as follows:

1. The processor decides to start a voltage change and asserts *VDD_Change*. Since this signal is generated within the core domain, it is synchronous to *CLK*.

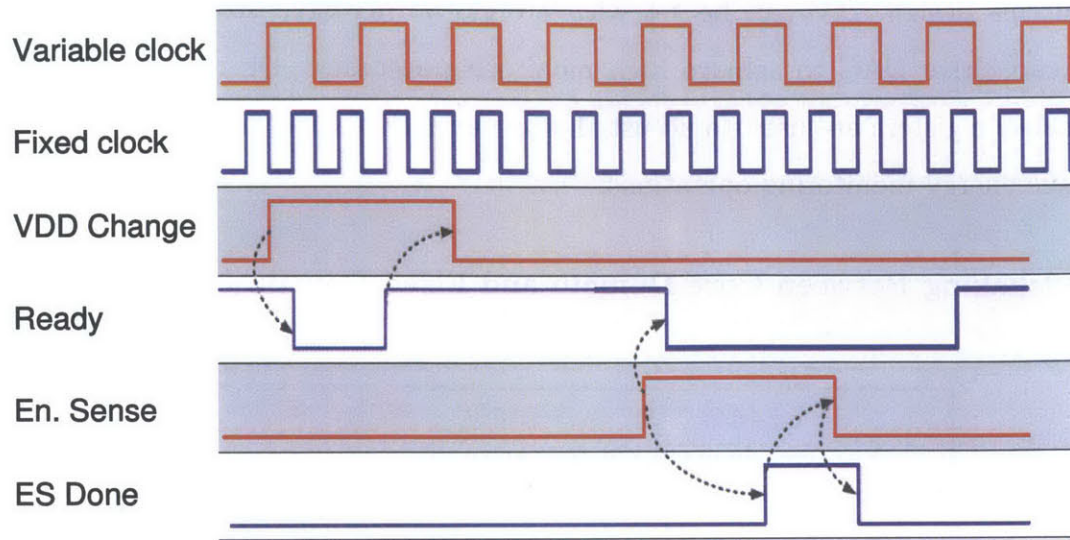


Figure 3-12: The asynchronous boundary signals between the fixed domain and the core domain.

2. The fixed domain state machine captures that *VDD_Change* is asserted, and deasserts the *Ready* which is synchronous to *FCLK*.
3. The new voltage level is reached, and the fixed domain state machine asserts *Ready*.
4. The processor captures that the *Ready* is asserted. Then, it deasserts the *VDD_Change* to finish the operation.
5. At this time, *VDD_Change* and *Ready* are at their initial values of '0' and '1' respectively. Thus, the system is ready for a new voltage change or an energy sensing cycle.

On the other hand, the operation of the handshaking for energy sensing cycle can be summarized as follows:

1. The processor decides to start an energy sensing cycle and asserts *Energy_Sense* which is synchronous to *CLK*.
2. The fixed domain state machine captures that *Energy_Sense* is asserted, and deasserts *Ready*. The energy sensing operation begins.

3. When energy sensing is done, the fixed domain asserts *ES_Done*.
4. The processor captures *ES_Done* is asserted. Then, it deasserts the *Energy_Sense*.
5. The fixed domain captures *ES_Sense* is deasserted. Then, it deasserts the *ES_Done*.
6. When the voltage is recovered to its original value *Ready* is asserted.
7. After this time, *ES_Sense*, *ES_Done* and *Ready* are at their initial values of '0', '0' and '1' respectively. Thus, the system is ready for a new voltage change or an energy sensing cycle.

Reference Voltages for Energy Sensing and Voltage Change Operations

Another challenge of the system was the requirement of multiple reference voltages for energy sensing and voltage change operations. First of all, for the voltage change operation and for creating the *Done* signal, two reference voltages are necessary. One of them is the high limit, V_{UP} and the other one is the low limit, V_{LOW} . Similarly, for energy sensing operation another low limit is required. Since in this system there are two domains, this means that six DAC blocks would be needed to create those six separate voltages. It is important to realize that in this design, either a voltage change or an energy sensing operation is performed at a given time. Therefore, DAC that creates the V_{LOW} can be shared between these two operations. This way, both operations can be performed using a total of four DACs.

The reference voltage creation during voltage change and an energy sensing operations are illustrated in Figure 3-13. To create the V_{UP} and V_{LOW} , 5 bit capacitive divider type DACs are used. The V_{UP} DAC is only required for the voltage change operation and its 5 bit input is delivered by the core. This 5 bit digital input is the sum of the digital representation of the new V_{DD} plus a differential. The digital value of the differential can change by 50 mV step sizes. On the other hand, the V_{LOW} DAC either creates the low limit for voltage change or energy sensing operation. For

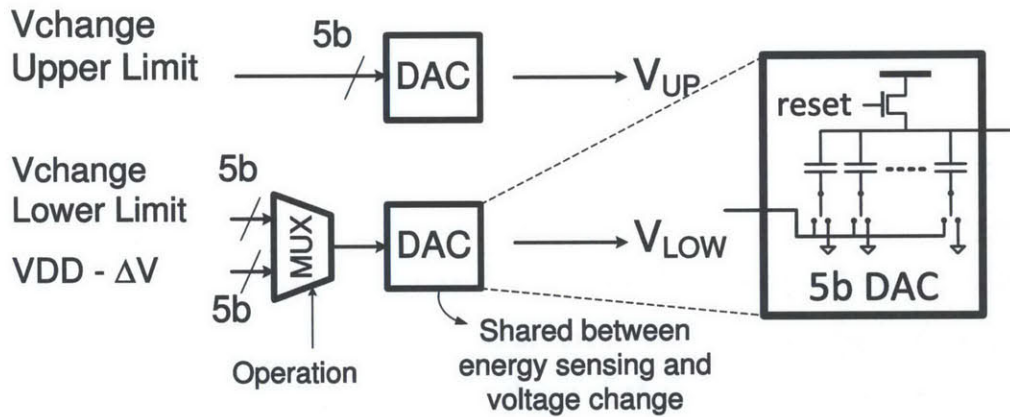
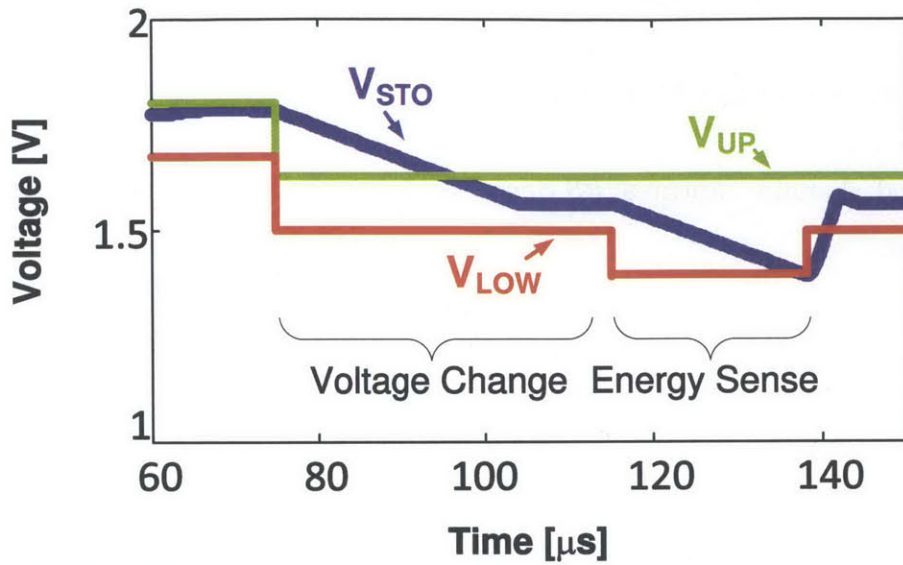


Figure 3-13: The reference voltage levels for energy sensing and voltage change operations.

voltage change, a similar 5 bit input is given by the core, which is V_{DD} minus a differential. For energy sensing, the input is $V_{DD} - \Delta V$.

3.4.3 Results for the DC-DC Converters with Embedded Energy Monitors

Efficiency Curves

Designing a highly efficient DC-DC converter is very important since power loss due to conversion can limit the benefits that are achieved by voltage scaling. For this

reason, the target in this design is to achieve over 90% efficiency. To achieve this efficiency, a buck converter with high efficiency is used. In addition, T_{PMOS} and T_{NMOS} pulse widths are designed to be reconfigurable.

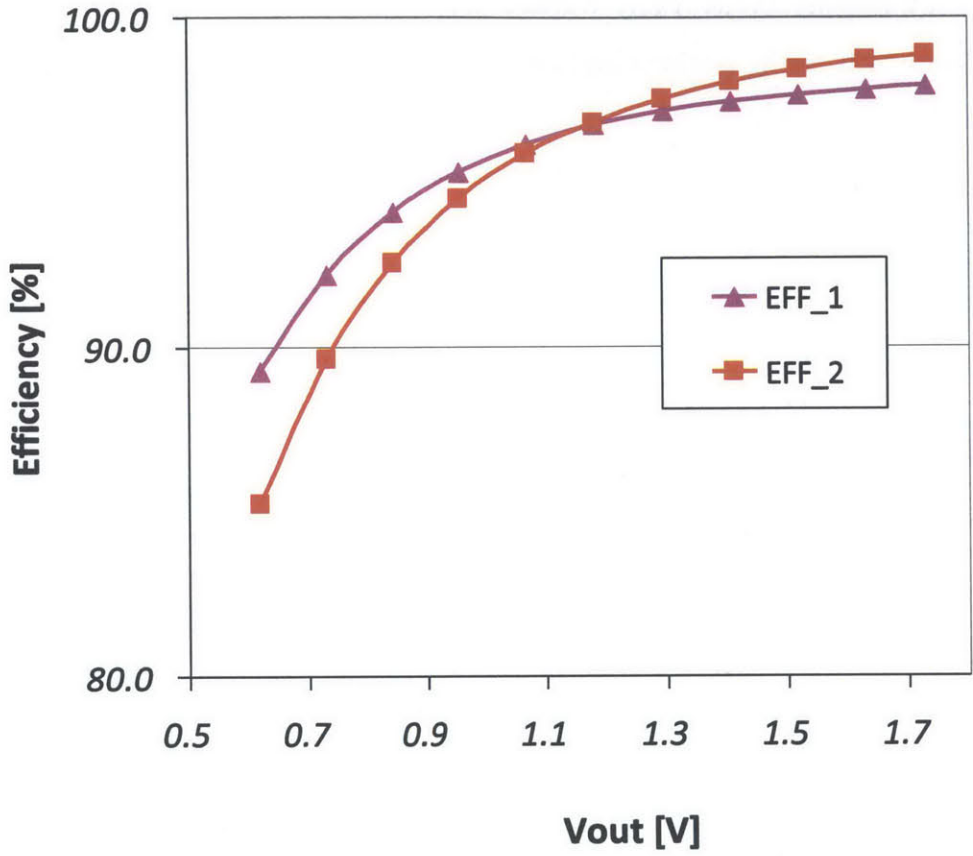


Figure 3-14: The superior efficiency curve can be selected for better efficiency as shown in the simulation result.

The PFM modulated buck converter efficiency is calculated as the *EFF1* curve shown in Figure 3-14. As mentioned before, this efficiency curve is achieved by fixing T_{NMOS} and modifying T_{PMOS} with V_{out} . Changing those two pulse widths together by the same ratio would result into an efficiency curve with a different slope. *EFF2* curve shows the efficiency curve that is achieved when both NMOS and PMOS pulse widths are halved in the design. Those two curves have different slopes due to the fact that *EFF1* curve depends on conduction losses more and switching losses less compared to *EFF2*. The actual slopes depend on different design choices such as the process selection, transistor sizes and parasitics. Designers would prefer to work on

the more advantageous curve at all times.

Although not implemented in this design, in order to determine which curve has a better efficiency during runtime, energy monitoring circuit operation can be used. After an energy sensing cycle, V_{STO} voltage drops by a fixed ΔV . Core can select different pulse width configurations and count the number of cycles it takes for the DC-DC converter to recharge the storage capacitor by ΔV using a fixed frequency. Then by comparing the number of cycles, it can determine which curve is more efficient.

For instance, for the design say that first configuration uses the original pulse widths and finishes in $T1$ time. Similarly second configuration uses two times shorter pulse widths and finishes in $T2$ time. By looking at the relationship between those two cycle counts, core can determine which curve is more efficient during runtime. If those two curves are equally efficient, $2 \times T1 = T2$. If $2 \times T1 < T2$, *EFF1* curve is more efficient and if $2 \times T1 > T2$, *EFF2* curve is more efficient. Cycle counting can be done multiple cycles to increase the precision in calculation of $T1$ and $T2$.

Oscilloscope Output

Figure 3-15 shows the oscilloscope output of the system while performing energy monitoring operations and voltage changes. In this example, first, a voltage change operation from 1.8 V to 1.7 V is performed. Then, it is followed by an energy monitoring period with a ΔV of 100mV. Then, the system can decide to scale the operating voltage down or up depending on its decision. This example illustrates the case where the system would decide to decrease the voltage to 1.6 V and perform a second energy monitoring period. This time, it chooses a 50 mV for the ΔV .

Based on the system runtime dynamics and targets, the system can perform multiple voltage change and energy sensing operations. Then, it can reconfigure itself for better power and performance optimization based on on-fly system targets that are set in the software.

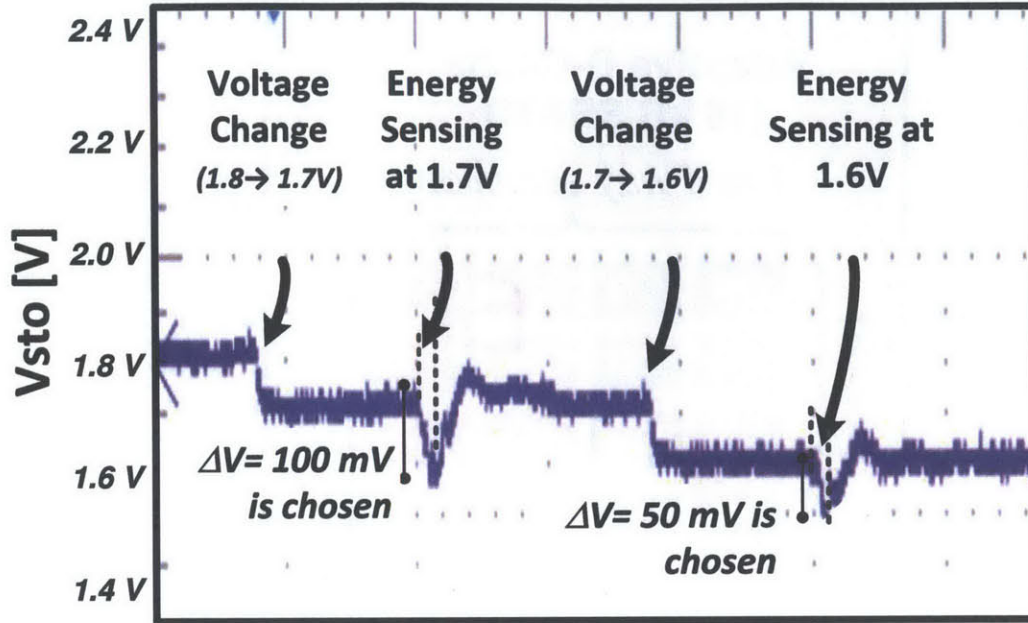


Figure 3-15: The oscilloscope output of the system while performing energy sensing and voltage change operations.

3.5 Reconfigurable SRAM Design

Embedded memories are one of the fundamental building blocks of the self-aware processor since they account for almost 60% of the total active area of the system. In the self-aware processor system, both the d-cache and i-cache are designed using SRAM memories to provide a high-density solution. Compared to implementing them with registers, SRAM design results into more than $4\times$ smaller area.

In the self-aware processor, both i-cache and d-cache are designed to be 16 kB and configured as sixteen 1 kB sub-blocks. The 1 kB memory sub-blocks are configured as 256×32 bits. Furthermore, for each 1 kB data memory, there is a 128 B tag memory block. The tag memories are also designed using SRAM memories and each 128 B memory is configured as 32×32 bits.

To enable operation across a large voltage range from 1.8 V down to 0.6 V, custom SRAM memories are designed using 8T bit-cells. The 8T bit-cell decouples the read operation from write operation and does not possess the read problem of its 6T alternative. To improve write-ability at low-voltages, peripheral row-drivers boost the word-line voltage up. 8T bit-cell based SRAM design will be explained in more

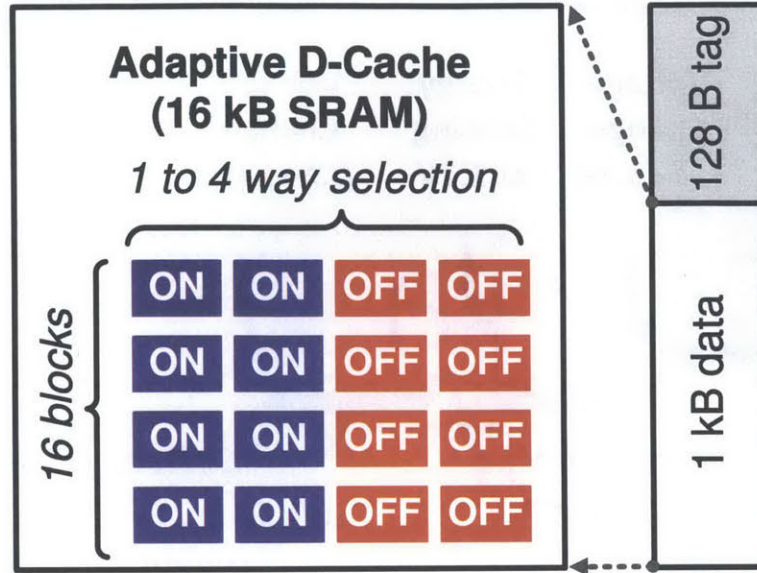


Figure 3-16: Adaptive d-cache structure.

detail in Chapter 4.

The d-cache memory is designed with dynamic associativity and size scalability. During runtime, the d-cache can be configured to be 1- to 4-way set associative and the size of each set can be configured to change from 1 kB to 4 kB. A configuration with a 2-way set associative d-cache with 4 kB memories per set is illustrated in Figure 3-16.

Resizing of the memories is possible thanks to the *Enable* input of the custom design SRAM memories (not shown on the figure). When this input is disabled, the memories are clock gated internally. In the example shown in Figure 3-16, a total of 8 sub-blocks of the d-cache are clock gated so they appear *OFF* and consume a significantly lower power. Clock gating 15 sub-blocks of the d-cache results into lowering the current consumption of the overall system by around 32% compared to keeping all 16 blocks *ON*.

Clock gating is a popular power reduction technique used in synchronous circuits for reducing dynamic power dissipation. On the other hand, power gating, which is an alternative technique to clock gating, reduces the power consumption by shutting off some parts of the blocks that are not in use. Power gating is generally a more effective method since it results into better energy savings. However, for SRAM

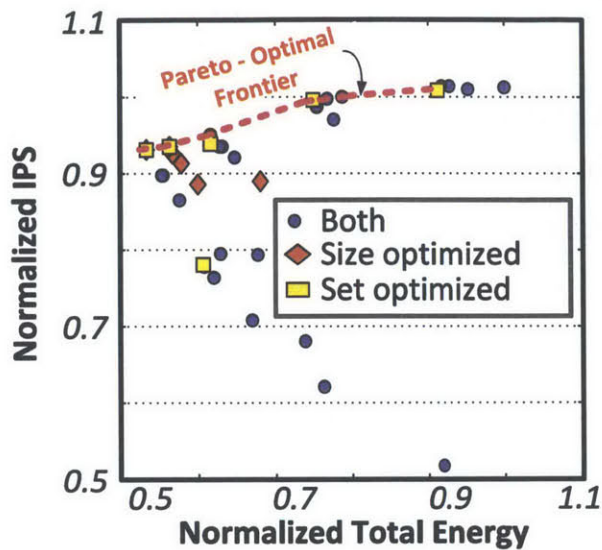
applications, current power gating methods can cause the bit-cells to lose their states; and, preserving the stored data in the caches during the *OFF* state might be desirable for some applications.

Lowering the supply voltage of the SRAM down to its data retention voltage (DRV) can be used as an alternative to power gating. This lowers the leakage consumption of the array due to DIBL while preserving the stored data. However, for technologies such as $0.18\mu\text{m}$, the leakage energy consumption might not be as significant as the dynamic energy consumption. Therefore, in this design, clock gating is implemented for turning off the unused SRAM sub-blocks. For more scaled technologies, where leakage is a significant portion of the total energy consumption, reducing the voltage to DRV can result into better savings.

3.5.1 Effect of Set Associativity and Size of L1-Cache

Associativity is a trade-off. Increasing associativity results into fewer cache misses so the core would waste less time to read from the off-chip memory. On the other hand, it needs to check more addresses which increases the complexity of the design and potentially the power consumption also increase. Similarly, cache size also has a strong effect on performance. In larger caches, there is a less chance that there will be a conflict. Again, this means that the miss rate would decrease when the memory size is increased.

The effect of changing the size and associativity of the L1-cache has an important effect on the total energy consumption of the system and the total time it takes for an instruction to finish. To illustrate the problems of the *closed* complex systems, we present the following experiment. It is performed in Graphite simulator [78] which is an open-source, parallel simulator for multi-core architectures. This simulator is developed in MIT CSAIL for multi-core processors. We run matrix transpose application on a single-core system with two possible adaptations: the total size of L2 cache (from 1 KB to 16 KB by powers of 2) or set-associativity (from 1 to 16 by powers of 2). For each combination of cache size and associativity, we measure the total time it takes to finish the application and the total energy consumed. The



App Name	Cache	
	Size	Set Assoc.
1d jacobi	8X	2-way
2d jacobi	2X	1-way
Matrix Transpose	1X	2-way
N-body	1X	2-way

Figure 3-17: Graphite simulation for changing cache size and associativity on a single core system. Different applications require different optimizations for minimum energy.

results are shown in Figure 3-17. The x-axis shows the total energy consumed and y-axis shows instruction per second (IPS). The blue circle points represent all tested configurations. The squares show configurations for a system that only considers cache set-associativity adaptations. Similarly the red diamonds represent points that only consider cache size adaptations. The best configurations are the ones that provide highest performance for lowest energy which is shown as the Pareto-optimal frontier. Notice that both squares and diamonds that appear to the right of the Pareto-optimal frontier are the points that closed systems would believe to be optimal, but in fact, are sub-optimal for the overall system.

Figure 3-17 also lists the optimum configuration of cache sizes and associativities for different applications when running on the same single-core system. This list shows the optimum configuration for the target of achieving the minimum total energy without a performance constraint. As it can be seen from the table, the optimum configuration is different for each application.

The self-aware processor system can dynamically adjust its cache size and as-

sociativity based on the current work load. A performance constraint can be set in the software and the system can minimize the energy with or without a specific performance constraint. As it can be seen from the figure, systems that enable both associativity and size adaptations can achieve up to 2× better IPS per energy. Therefore, reconfigurability of d-cache is an important technique to reconfigure the system under changing dynamics.

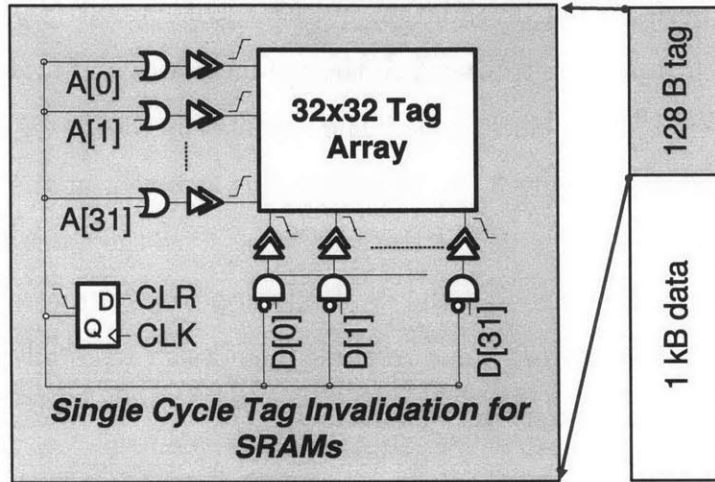


Figure 3-18: Single-cycle tag invalidation.

3.5.2 Single Cycle Tag Invalidation

For every data stored in a cache, there is also a tag part that needs to be stored as well. The data block contains the actual data fetched from the main memory. On the other hand, the tag contains a part of the address of the actual data. This is how the processor decides whether a data in the memory is cached or not. In most designs, data cache tag entry consists of two fields in the tag memories: address bits and valid bits. For the self-aware design, the most significant 22 bits out of 32 ($TAG[31:10]$) are the address tag and they contain the address of the data held in the cache line. The least significant 4 bits ($TAG[3:0]$) are the valid bits. When set, the corresponding sub-block of the cache line contains valid data. These bits are set when a sub-block is filled due to a successful cache miss.

After certain reconfigurations such as expending the size of the memory, the tag

memories need to be invalidated. In this configuration, the tag memories are 32 lines and there is a total of 16 blocks of tag memory. Therefore, if implemented without any special circuitry, a tag invalidation would require sequential writes to the tag bits of the memory which would increase the overhead of cache reconfiguration. Therefore, it is important to be able to perform tag invalidation as quickly as possible.

The self-aware processor tag memories are designed with a special tag invalidation circuitry to perform tag invalidation quickly. This circuitry is not standard in processors. A *CLR* input is introduced to the tag memories and it is used in order to clear all tag bits in a single clock cycle. The circuitry is shown for one block of the sixteen tag memories in Figure 3-18. When *CLR* is asserted, at the positive edge of the CLK, all the WL's are enabled at the same time since the address bits are OR'ed with the *CLR* input. Similarly, all the data bits appear as zeroes. Therefore, when asserted, the synchronous *CLR* input to these memories causes all 32 words of each block to be overwritten with zeroes simultaneously.

3.6 Measurement Results of the Self-Aware Processor

For proof-of-concept, the ideas explained in this chapter are implemented in a test-chip prototype using a $0.18\mu\text{m}$ CMOS technology. However, it would scale well with technology and can be applied in more scaled nodes. The specifications of the test-chip are given in Table 3.1 and the die photo is given in Figure 3-19. The total size of the test-chip is $6\text{mm} \times 6\text{mm}$. The core is implemented using the digital flow whereas the SRAM circuits (d-cache, i-cache and trace buffer) and the two DC-DC converters with embedded energy monitoring circuits are implemented using custom design circuits.

The two voltage domains of the system are powered by the two DC-DC converters. The DC-DC converters are placed at the corners of the die to minimize their distance to the pads. Similarly, in order to minimize the parasitic resistance and inductance

Table 3.1: Test-chip specifications for the self-aware processor SoC designed in 0.18 μm CMOS.

	Technology	180 nm LP CMOS
	Chip Size	6 mm x 6 mm
SRAM	Organization	8T with write-assists
	Data Memories	16 kB (256x32b blocks)
	Tag Memories	1 kB (32x32b blocks)
	Adaptation	Size 1-16kB, 1-4 Way
DC-DC	Resolution	5 bit
	Modulation	PFM
	I/O V_{DD}	1.8 V to 0.6-1.8 V
Energy Sensor	ΔV Step Size	50 mV
	Dynamic Power Overhead	< 0.1%
	Area Overhead	< 1%

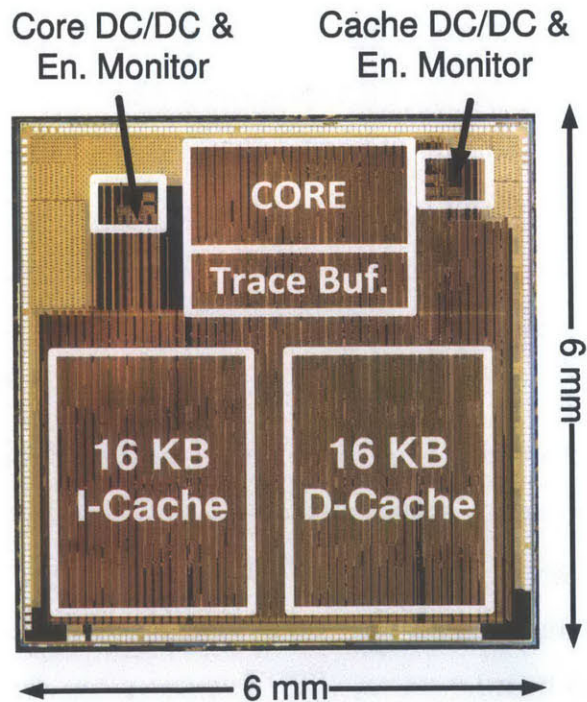


Figure 3-19: Die photo of the 0.18 μm self-aware processor system.

of the bond wires, multiple bond pads are connected to the same pin on the package. Two energy monitors are embedded into the DC-DC converters. The ΔV step size for energy sensing is 50mV.

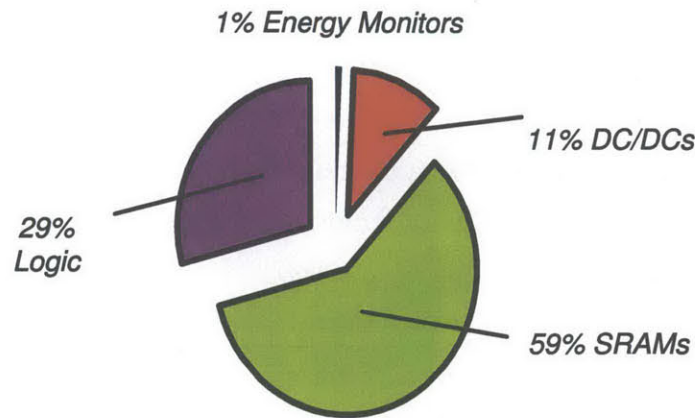


Figure 3-20: Area breakdown of the test-chip prototype.

The active area breakdown of the test-chip can be seen in Figure 3-20. The energy sensors are almost free since they bring a very small area ($< 1\%$) and power ($< 0.1\%$) overhead. On the other hand, the SRAMs account for almost 60% of the total area, which makes their design significant. The core design is around 35K gates and it is scattered in the 29% of the active area of the chip. The trace buffer, which is used to keep a copy of the executed instructions for debugging purposes, it is also implemented using the custom design SRAM memories.

The self-aware processor system is tested in a test setup that is illustrated in Figure 3-21. A daughtercard board is designed and the self-aware test-chip resides on it. On the daughtercard, there are also on-board DC-DC converters that generate the required I/O voltages of 3.3 V and 1.8 V from a 5 V supply. Daughtercard board also has status LEDs and buttons for testing. The filtering capacitors and inductors are also placed on this board.

The self-aware processor system outputs the AHB bus off-chip. An opal-kelly board, with a Xilinx Virtex-5 FPGA, talks to the AHB bus through the 240b expansion connectors. The I/O interaction between the main memory and the test-chip is performed through this connection. The 32MB DDR2 DRAM that resides on the opal-kelly board is used as the main memory. The *front-panel* creates a communica-

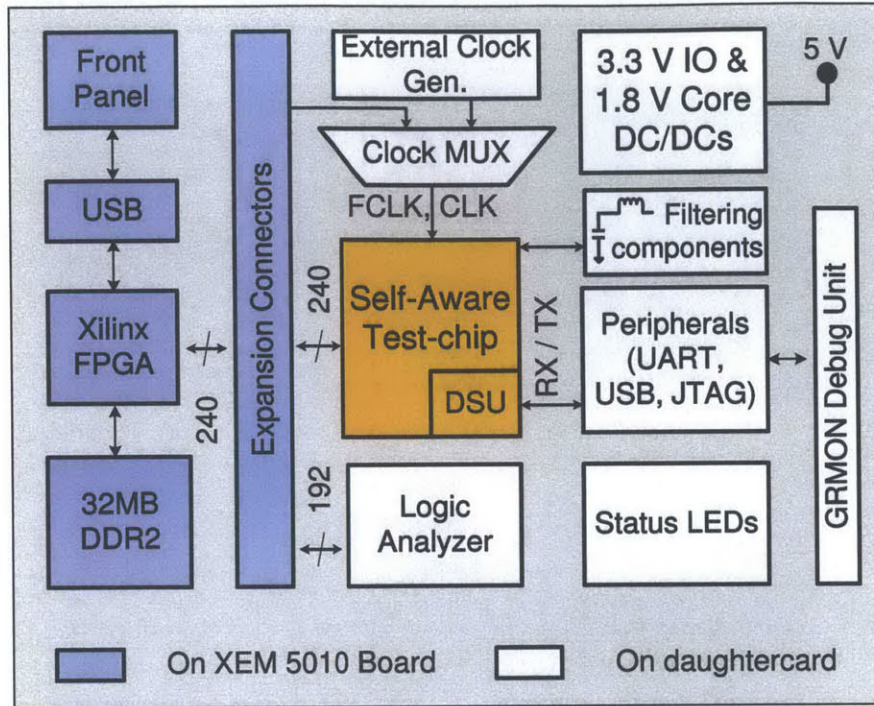


Figure 3-21: Block diagram of the test setup of the self-aware processor system.

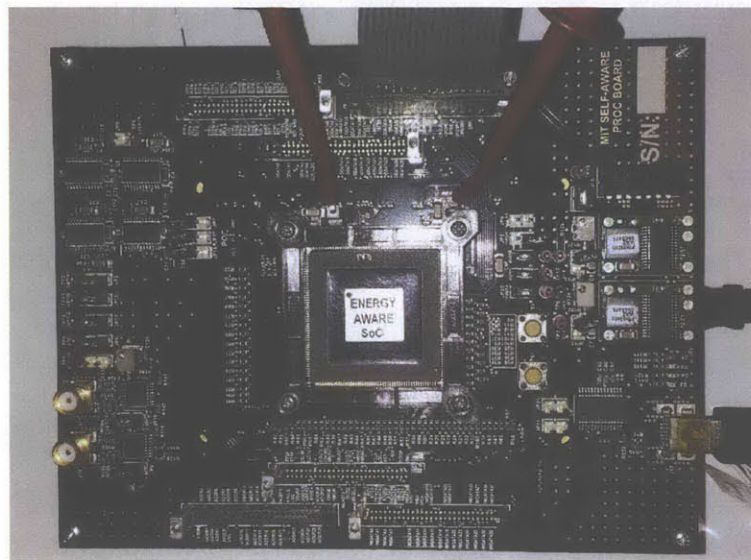
tion between the FPGA and the PC through a USB connection. The two clocks for the system can be created using either the FPGA or a high-speed clock generator.

GRMON, which is a debug monitor for the LEON3, is used extensively for testing purposes and to aid hardware and software debugging. Through the debug support interface, full access to all processor registers and caches is provided. It also enables step by step instruction execution. Additionally, the internal trace buffer monitors and stores the executed instructions, which can be later read out over the debug interface, DSU. This provides a non-intrusive debug environment on the real target hardware.

The picture of the board setup can be seen in Figure 3-22. The opal kelly board is designed to connect the daughtercard board through the expansion connectors. The daughtercard board resides on the top whereas the opal kelly board sits at the bottom. This way, the path delays for the AHB and CLK signals are kept as minimum.



SIDE VIEW



TOP VIEW

Figure 3-22: Picture of the test setup of the self-aware processor system.

3.6.1 Effect of DVFS and D-Cache Adaptation to Energy Consumption

Figure 3-23 (a) shows the measured energy consumption of the system running a matrix transpose benchmark, APP1 under the effect of voltage and frequency scaling. Total EOP scales from 3.85 nJ to 690 pJ with DVFS only. The system is measured to achieve a 65MHz clock frequency at 1.8 V and 22 MHz at 0.82 V.

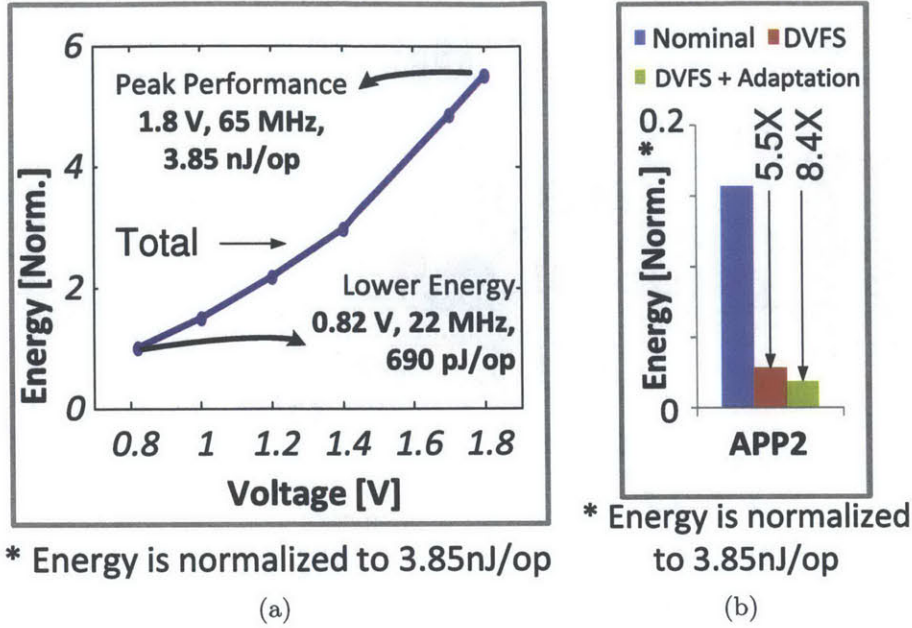


Figure 3-23: (a) The effect of DVFS on energy consumption, and (b) the effect of DVFS and cache adaptation.

The system is also capable of circuit reconfigurability on top of DVFS. For the same application, how much energy savings can be achieved under both DVFS and cache reconfiguration is measured. As it can be seen in Figure 3-23 (b), up to 8.4× lower energy is achieved through both DVFS and dynamic adjustments of d-cache size and associativity. The energy savings due to DVFS only is 5.5×. A 1.54× energy savings is achieved on top of DVFS with the help of cache reconfigurability compared to operating the system at 1.8 V and full memory.

3.6.2 Effect of D-Cache Adaptation for Different Applications

Figure 3-24 shows a scatter plot of the measured power and performance trade-offs of the self-aware system under cache configuration for two applications:

1. APP1 is a matrix transpose application for a 16×16 matrix size.
2. APP2 is the same application for a 32×32 matrix size.

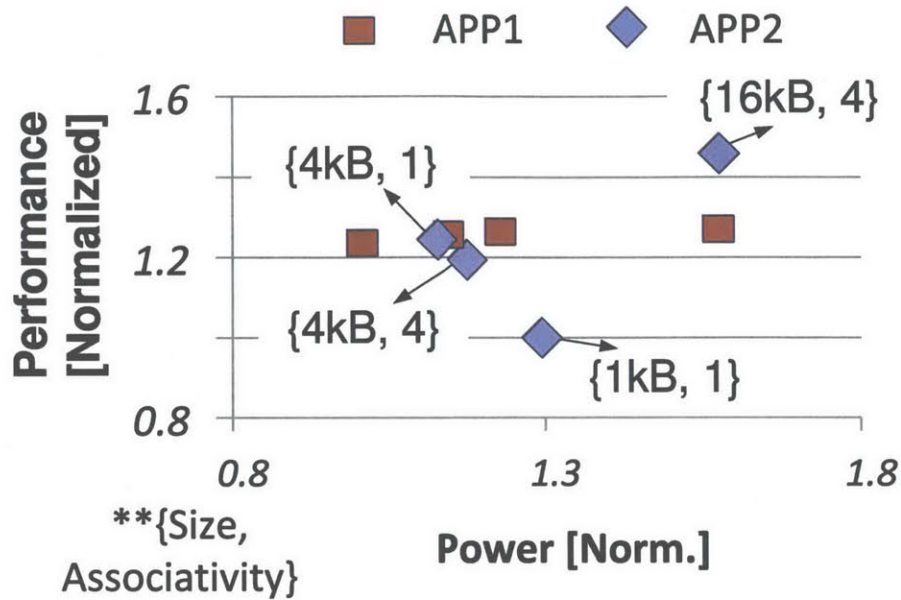


Figure 3-24: The effect of cache adaptation on energy consumption.

For ease of illustration, only four points per application are shown and each point represent four corner cases for the d-cache configuration for size and associativity. Each point represents an operation using a different cache configuration at the same voltage (1.8 V) and clock frequency (CLK=35 MHz). The y-axis is the normalized performance and it is proportional to the inverse of the time it takes for the application to finish using a specific configuration. The x-axis is the normalized total power consumption, which is the total power of the core, caches and the I/O. For an application, the preferred cache configurations are the ones that reach the same performance for minimum power consumption.

The optimum d-cache configuration is determined based on whether there is a performance constraint. For instance, for APP2; {size,set associativity} = {4kB,1} and {16kB,4} are on the pareto optimal frontier. Depending on the performance limit, one of those two points can be selected as the optimum. The other two configurations result into a smaller performance for higher power consumption compared to {4kB,1} configuration.

APP1 is not cache-bound since its working set can fit into the smallest cache configuration. This is the configuration of $\{\text{size, set associativity}\} = \{1\text{kB}, 1\}$. Therefore, for APP1, increasing the cache size increases the power consumption without any significant improvement in performance. On the other hand, APP2 has a larger working set which cannot fit into the smallest cache size but would be able to fit into an intermediate cache size. Therefore, its performance increases by increasing the cache size. Contrary to intuition, the smallest cache configuration of $\{1\text{kB}, 1\}$ does not result into the minimum power consumption for this application. This is due to the fact that it results in an increase in cache misses and total power consumption rises due to the larger I/O power. Therefore, if minimum energy consumption is desired for APP2, the system would choose to work with an intermediate memory configuration of $\{4\text{kB}, 1\}$.

3.6.3 Comparison to a Conventional Race-to-Idle System

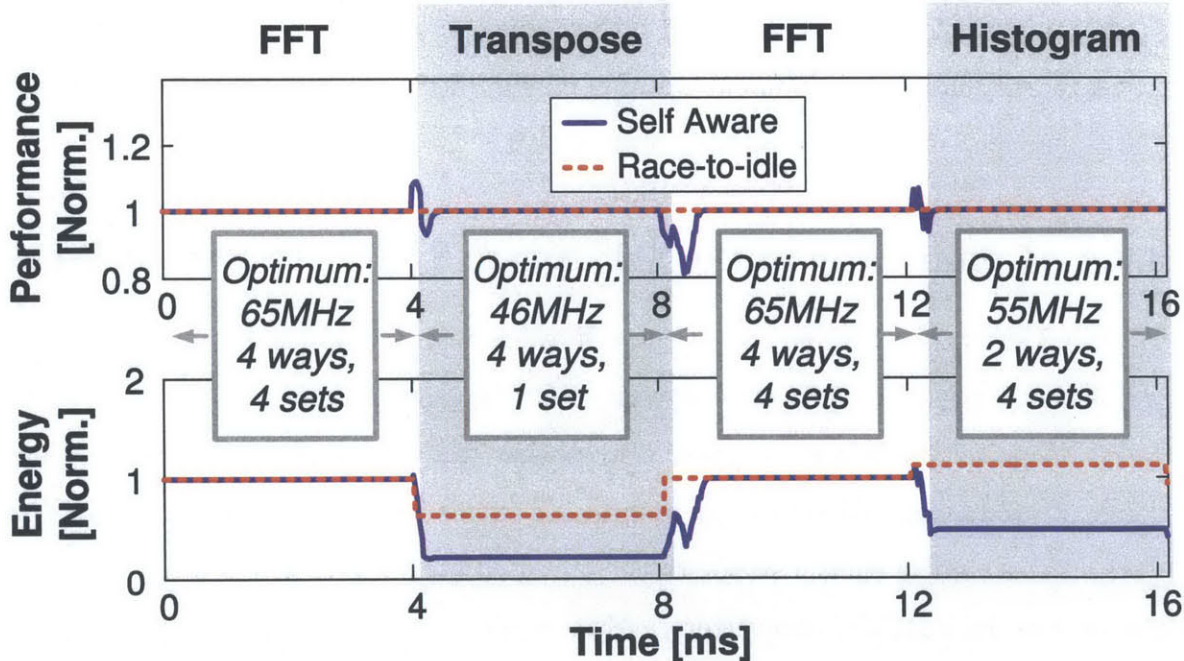


Figure 3-25: System simulation running four phases of a multi-media application using 1-self-aware adaptation, 2-static configuration with race-to-idle operation.

Figure 3-25 shows a system simulation where SEEC is optimizing the system

for a multimedia application with four distinct phases: FFT, transpose, FFT, and histogram. This sequence can be a basic image processing pipeline that applies a frequency domain filter and then computes a histogram.

For a standard system that does race to idle and cannot reconfigure itself, the system needs to be designed with enough resources to meet the performance needs of the most demanding phase. For this example, the most demanding phase is FFT and it would require 4ms to finish using all the resources. If other phases do not need those resources (full cache size, associativity and highest voltage case), they complete before their deadline and then idle. Therefore, the race-to-idle system meets the performance constraint for every four phase. This case is shown with the dotted red line.

On the other hand, the other line is the self-aware system that exposes adaptations to hardware. The four distinct phases are shown on the figure with the configuration SEEC chooses for that phase. While meeting the same performance goal, the self-aware design achieves an almost $2\times$ reduction in energy compared to a design with a static configuration and conventional race-to-idle operation. For a given performance target, each phase has a different optimal configuration (ways, sets and clockspeed) and the proposed system is capable of finding it.

3.6.4 Comparison of the Self-Aware System to Recent Work

Table 3.2 shows the comparison of the self-aware test-chip to the recent work. Here [22] is the Sandybridge processor with its PCU blocks that generate power numbers based on power modeling. [72] is the work that was explained in Chapter 2 which demonstrated an embedded energy monitoring circuit for an SRAM application. This work used a similar energy sensor to the one used here, however it did not integrate the sensors into DC-DC converters. Also, it did not do any adaptation using the energy readings. The third work is [33] shows a system with a unique energy sensor. However, it also does not use the information for controlling the system. The work in [76] also illustrates a self-aware system using relative energy comparisons. It does an excellent job of optimizing the system for the minimum energy point using the energy

Table 3.2: Comparison of the self-aware test-chip to previous work.

Work	[22]	[72]	[33]	[76]	This Work [73]
Adaptation	Core,L3 cache, DVFS	DVFS	no. adapt.	DVFS	DVFS, d-cache
V_{DD} [V]	0.7-1.15	0.37-1.2	3.3	0.25-0.7	0.6-1.8
En. Monitor Accuracy	model based	10%	20%	NA	10%
En. Monitor Area Overh.	NA	16%	NA	21%	1%
En. Monitor Area Overh.	thermal sensor	capacitor (off-chip)	PFM DC-DC	capacitor (shared)	capacitor (shared)
En. Savings (due to DVFS)	NA	4.8×	no DVFS	4.1×	5.5×
En. Savings (due to adapt.)	NA	no adaptation	no adaptation	no adaptation	1.5×

readings. However, the system cannot work with changing power and performance constraints since it always tries to minimize the energy consumption without any performance constraint. Also, the energy sensors do not strictly limit the voltage ripple on the power supply.

The work that was explained in this chapter is summarized in the last column [73]. It uses multiple energy sensor readings for reconfiguring system components for better power and performance optimizations. The energy sensors used in this work can generate absolute EOP readings by actually measuring the real energy consumption.

3.7 Summary and Conclusions

Modern processor systems are getting more and more complex at every node. Similarly, they must balance multiple and often competing design goals such as maximizing performance while minimizing energy. Furthermore, they have to work optimally under dynamic operating conditions such as temperature and voltage fluctuations, process variations, aging, and with a wide variety of applications with different phases.

Self-awareness is a popular field of research in computing that considers systems and architectures that gather and maintain information about the current state and environment, reason about the behavior and adapt themselves if necessary

Adaptive systems which leverage power management engines based on power models are used to improve power and performance optimizations. However, power models cannot fully represent the actual profile of a complex processor system. We showed an absolute energy monitoring circuit for an SRAM application in Chapter 2, but integrating them into microprocessors would be beneficial. Furthermore, enabling the hardware reconfigurability as a software knob would result into the software to control the adaptation under changing power and performance constraints. Furthermore, additional benefits can be obtained by integrating those energy monitors within DC-DC converters.

This chapter presents a self-aware processor system with two energy monitoring circuits that can measure actual energy consumption on the fly. The monitors are embedded into DC-DC converters and do not require any extra off-chip components since the off-chip filtering capacitor is shared. The design is based on a LEON3 single-core processor. Efficient power conversion for the two power domains of the system is provided by two on-chip DC-DC converters that deliver variable load voltages from 0.6V to 1.8V. One of the DC-DC converters powers the core and i-cache while the other powers the d-cache. Two energy monitors allow the system to distinguish between energy spent on computational operations versus energy spent on data storage. Performance counters are included to track dynamic performance changes. The instruction and data caches are constructed from custom-design SRAMs. Low-voltage SRAM operation is made possible through the use of 8T bit-cells and write-assists. The d-caches are designed to be re-configurable in associativity and size to adapt to compute- versus cache-bound phases of applications. Cache configuration is performed in < 3 clock cycles including tag invalidation.

These hardware features are open to software and a software self-aware computation engine (SEEC) can dynamically adapt the processor to meet performance and energy goals. Measurement results show that up to $8.4\times$ energy savings can be

achieved with DVFS and self-adaptation.

Chapter 4

Low-Voltage SRAM Design

On-die cache size is increasing rapidly over the years (Figure 1-3) and it is projected that this trend will continue in the future. Therefore, SRAMs are one of the fundamental building blocks of today's systems, and they account for a significant area and energy consumption. Thus, designing energy-efficient and high-density SRAM circuits has been and will continue to be an important research topic.

One popular way to increase energy-efficiency of the SRAM circuits is using voltage scaling. Figure 4-1 shows the bit-cell and V_{DD} trend of SRAMs from major semiconductor companies [3]. As it can be seen from the figure, in today's scaled technology nodes, the minimum operating voltage of the SRAM circuits is only 100mV lower than the nominal voltage. Therefore, operating SRAMs at lower voltages not only increase SRAM energy efficiency, but also enables them to be powered from the same supply voltage of the logic circuit. This is very desirable since the system level design becomes simpler.

In this chapter, two different SRAM designs will be explained and both of them target low-voltage operation. Section 4.1 talks about SRAM design that is based on 8T bit-cell topology. This section also examines peripheral assist circuits to make 8T bit-cell based SRAM circuits work at low-voltages. Furthermore a SA offset compensation technique using body-biasing will be introduced. A prototype test-chip was shown in Chapter 2 that demonstrated an energy monitoring circuit, and this chapter will focus on the SRAM circuits inside it. Section 4.2 will talk about a

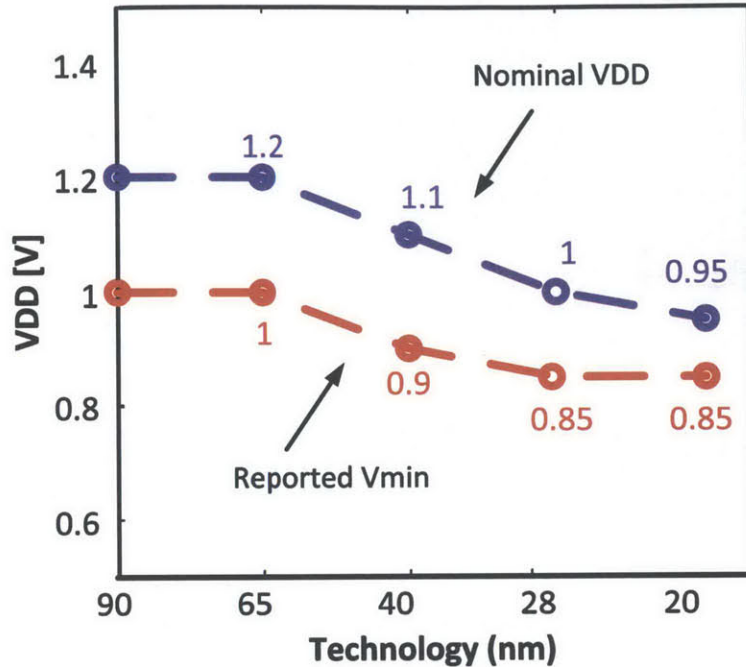


Figure 4-1: Bit-cell voltage scaling over technology nodes [3].

6T bit-cell based memory. A prototype test-chip is designed using a 28nm FD-SOI technology. This design is measured to be operational down to 0.43V using industry sized 6T bit-cells. In this section, the assist circuits that are used to achieve low-voltage operation for this 6T bit-cell based SRAM will be examined.

4.1 8T Bit-cell Based SRAM Design

8T bit-cell based SRAM design is an alternative to 6T and generally it can achieve a lower- V_{DD} operation. Thus, in this thesis, two test-chips are implemented using 8T bit-cell based SRAM circuits:

1. An embedded energy monitoring circuit was illustrated for an SRAM application in Chapter 2. Those SRAMs are designed using 8T bit-cells using a 65nm LP CMOS process and they are measured to be voltage scalable from 1.2V down to 0.37V.
2. In Chapter 3, we showed a self-aware processor which was designed using a $0.18\mu\text{m}$ technology. The d-cache, i-cache and trace buffers of this system were

implemented using 8T bit-cells and they are voltage scalable from 1.8V down to 0.6V.

For the rest of this section, we will focus on the former test-chip prototype which was designed using the 65nm CMOS. The latter design uses similar design ideas and therefore it will not be covered separately in this thesis.

4.1.1 8T and 6T bit-cell Topologies

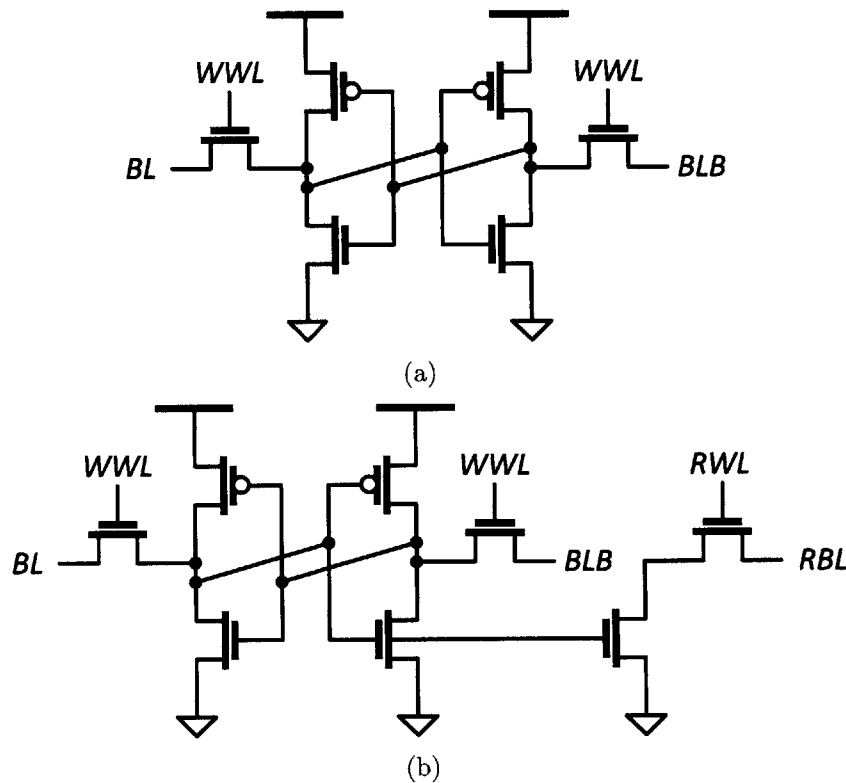


Figure 4-2: Schematics of (a) a 6T bit-cell, and (b) an 8T bit-cell.

The schematics of the conventional 6T bit-cell and non-conventional 8T bit-cell are shown in Figure 4-2. The 6T bit-cell has been used as the fundamental building block for the SRAMs for many years due to its compact layout and design. However, as explained in Chapter 1, its operation depends on the relative strengths of its transistors. Therefore, it is a ratioed circuit and its operation suffers from read, write and retention related problems at low voltages.

An alternative 8T bit-cell is proposed in [79]. This bit-cell has two extra transistors, constituting the read buffer. In this bit-cell, a write operation is performed through write word-line (WWL), BL and BLB ports whereas single-ended read is exercised through read word-line (RWL) and read bit-line (RBL) ports. Therefore, the read and write operations are decoupled from each other and read-upset problem of the 6T bit-cell does not exist in 8T bit-cell. Although 8T bit-cell is more promising in terms of voltage-scaling, it still has some challenges.

First of all, although 8T bit-cell does not have the read-upset problem, it still experiences the same write and retention problems of the 6T bit-cell. One possible solution is that, since read and write operation are decoupled, the 6T part of the 8T bit-cell can be designed to favor writes. However, transistor sizing brings diminishing returns to write-ability since variation depends on $1/\sqrt{W \times L}$ [80]. Therefore, 8T bit-cell can benefit from extra write-assist techniques to lower its V_{min} .

Secondly, it has around 40% larger area compared to a 6T bit-cell [56, 79]. For designs where density is the highest priority, this area overhead might be unacceptable.

Thirdly, it does not work with column interleaving in the conventional way. This is due to the fact that during column interleaving, the half selected bit-cells are under read stress which makes the 8T bit-cell lose its advantage over 6T bit-cell. Therefore, the periphery circuits for 8T bit-cells need to be larger since each column requires a dedicated read and write circuitry. [61] proposed a column-interleaved SRAM design using 8T bit-cells. However, it brings around 10% area overhead.

Moreover, without column interleaving, 8T bit-cells are more susceptible to soft errors. This is particularly important for cases where multi-cell upsets (MCU)s cause multiple fails in the same word and may not be corrected using a single error correcting (SEC) ECC. Alternatively, a more complicated ECC scheme can be used for the 8T bit-cell to overcome this problem. For instance, a 6T bit-cell based memory with 2 to 1 column interleaving can use a SEC ECC to be resilient up to two bit MCU. For a 128b word, a SEC ECC requires 8 parity bits. On the other hand, for the same resiliency, an 8T bit-cell would require a double error correcting (DEC) ECC which

would require 16 parity bits [81].

4.1.2 Error Map Comparison for 6T vs. 8T

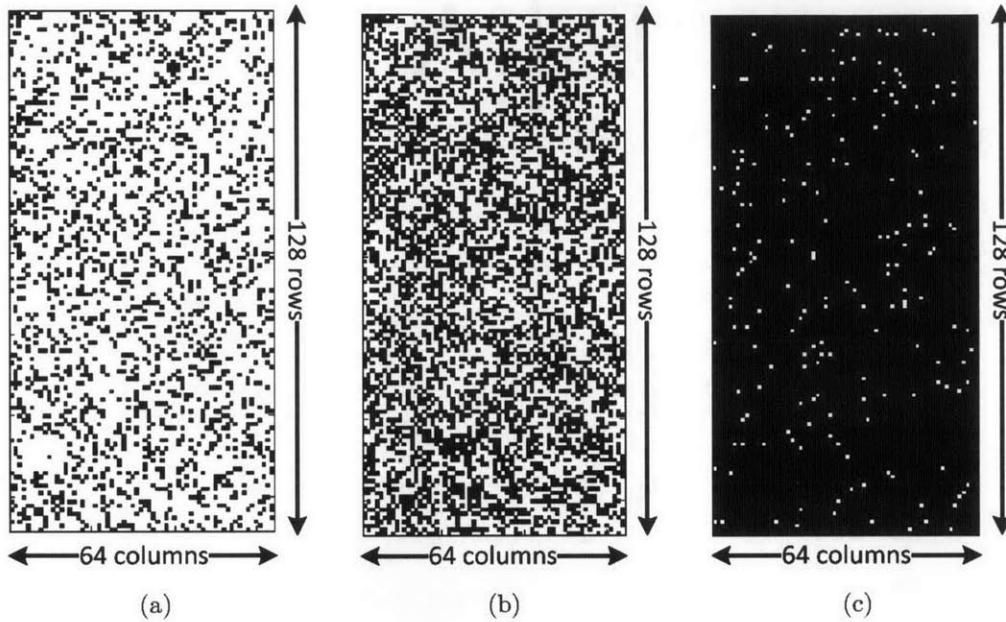


Figure 4-3: Error map simulation of a 64×128 SRAM block designed using a (a) 6T bit-cell, (b) 8T bit-cell without assist, and (c) 8T bit-cell with write-assist. $V_{DD}=250\text{mV}$. White: erroneous bit; Black: correct bit.

Figure 4-3 shows the error map of an SRAM block with 128 rows and 64 columns at 250mV. For this simulation, we used a predictive 22nm technology [82]. The white dots represent erroneous bit-cells whereas black dots represent correctly operating bits. The errors can be due to any type of read, write or retention.

- Figure 4-3 (a) shows the error map for a 6T bit-cell based memory.
- Figure 4-3 (b) shows it for an 8T bit-cell based memory without any assist circuits.
- Figure 4-3 (c) shows the error map for an 8T bit-cell based memory using write-assists.

As it can be seen from the figure, 8T memory has fewer errors compared to its 6T counterpart since it does not have read-upset problem. Similarly, the 8T bit-cell based memory with write-assists performs better compared to the 8T memory without assists since it has better write-ability. Therefore, in the 8T bit-cell based memories, write-assist techniques are used to make them work at low voltages.

4.1.3 Write-Assist Circuit

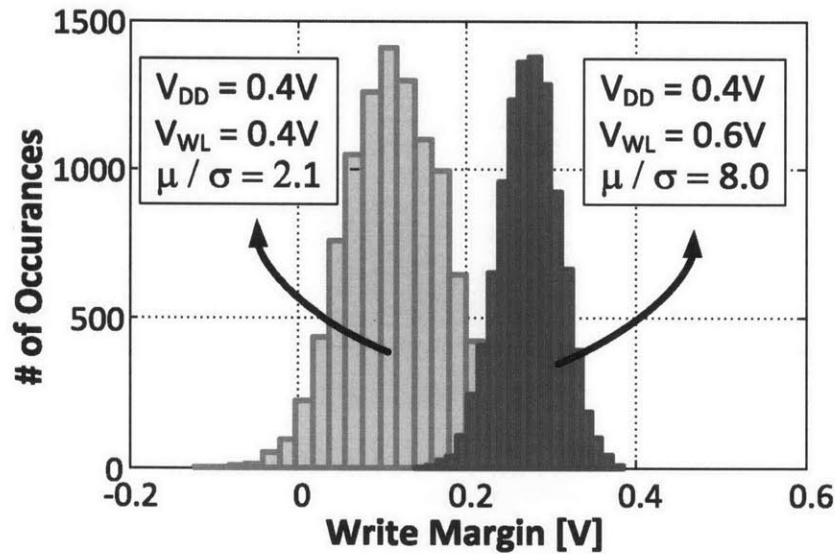


Figure 4-4: For a 65nm technology, 10K MC simulation shows that 200mV higher V_{WWL} improves write failure σ/μ by $3.8\times$ at 400mV.

8T bit-cell experiences write failures at ultra-low voltages. Therefore, to provide a robust write operation at low voltages, a write-assist scheme can be necessary. Different techniques have been proposed for improving write-ability in the literature. Some of those techniques can be summarized as lowering the bit-cell supply voltage, WL boosting, using a negative BL and body-biasing as mentioned in Section 1.3. A summary of different assist techniques can be seen in Appendix B.

In this design, this issue is addressed by using an increased supply voltage for WWL drivers. This technique cannot be used for a 6T bit-cell the same way since it would degrade the read margins of the half-selected bit-cells on the same row.

However, for the 8T bit-cell based memory, the WWL is a separate signal from the RWL so this technique can be implemented.

Increasing WWL voltage by 200 mV over the bit-cell array voltage results into $4\times$ improvement in write margin μ / σ at 400 mV. This phenomenon can be observed in Figure 4-4. In this analysis, 10k MC simulations are performed. A negative margin indicates failing memory bit-cells. As it can be seen from the figure, many bit-cells fail to work at 0.4 V if no write-assist circuit is used. On the other hand, using $V_{WWL}=600$ mV results into a $\mu / \sigma=8$ which indicates that the memory is not expected to have write failures at 400 mV. Figure 4-5 shows how the write-assist is implemented in this design. As it can be seen, the WWL voltage is boosted up by a level converter near the bit-cell array.

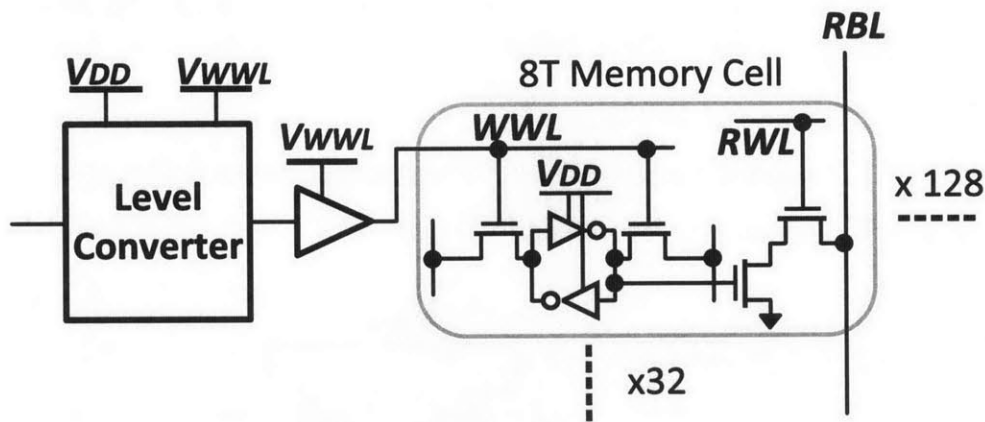


Figure 4-5: 8T bit-cell based design with write-assist enables operation from 1.2V down to 0.37V.

An important design decision related to this dual supply assist (V_{DD} and V_{WWL}) is how to partition the voltage domains. The simplest approach would be to supply all the circuits in the WWL path at the high voltage as shown in Figure 4-6(a). The second option would be to place the level converters internal to memory at the boundary between the bit-cell array and the row drivers Figure 4-6(b). The first option would result into a smaller area since it involves the lowest number of level shifters. However, this would result into a higher power consumption since a significant portion of the memory periphery would be at the high voltage domain. Using the level shifters inside the memory, the power consumption would be reduced

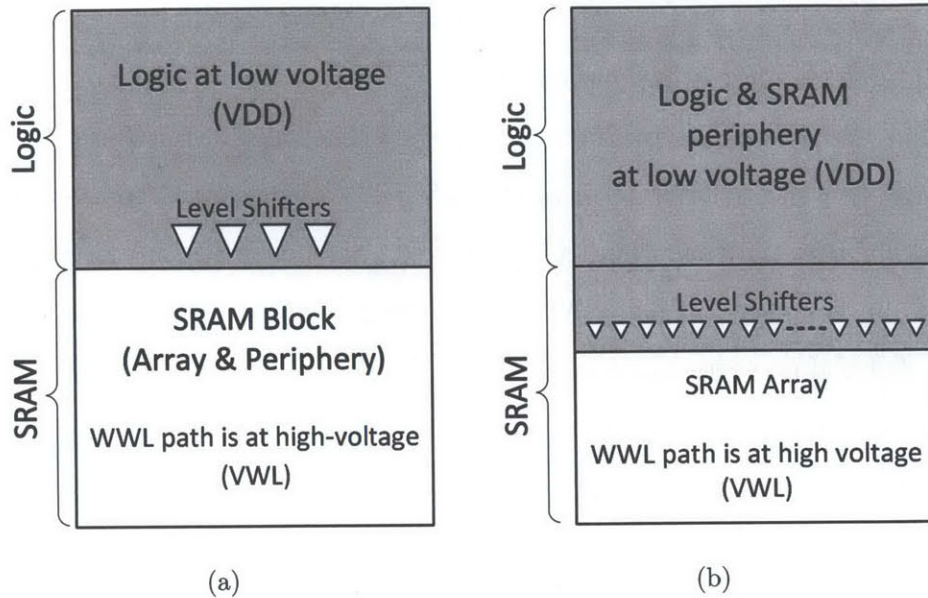


Figure 4-6: Two options for partitioning the voltage domains in dual supply write-assist technique (a) Level shifters between boundary and logic, (b) Level shifters internal to SRAM.

by around 20%. In this design, we chose to implement the second option as can be seen in Figure 4-6(b) since it favors energy for area.

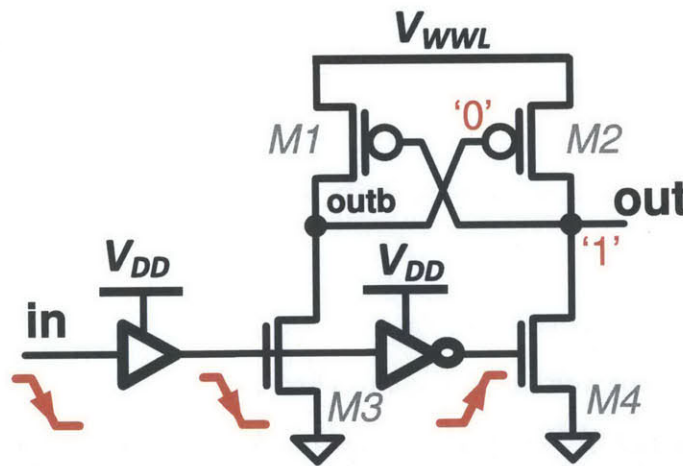


Figure 4-7: The schematics of the DCVSL type level shifter used in this design under level shifting operation.

As level shifters, differential cascode voltage switch logic (DCVSL) type is chosen. The schematics of this level converter design is given in Figure 4-7. This level converter does not dissipate any static power consumption. However, as it is well

known, it is a ratioed circuit and correct operation depends on the relative strength of the transistors. For instance, if the output is '1' and the input transitions from high to low, M_4 turns ON. However, at this time, M_2 is also ON and it holds *out* high. Therefore, for correct operation, M_4 needs to be able to overpower M_2 . The sizing of the transistors depends on the highest voltage difference between V_{WWL} and V_{DD} . In this design, the width of the NMOS transistors are chosen to be $16\times$ larger which guarantees correct operation under the effect of variations.

Since in this design, write-assists are used in the periphery, the access transistor sizings are kept as minimum. On the other hand, the read buffer transistors are chosen larger to improve the read current.

4.1.4 Memory Organization of the 128kb SRAM

Figure 4-8(a) shows the array architecture of the 128 kb SRAM which is structured as four 32 kb blocks with 256 rows and 128 columns. All four blocks have their dedicated row circuits, decoder and column circuits. The SRAM blocks use a two stage (or hierarchical) sensing which will be explained in more detail in the next section. In terms of organization, every block consists of 8 sub-blocks, dividing the 256 rows of the block into 32 rows per sub-block. This way, 32 bit-cells on a sub-block share a common bit-line shown as *FBL*. On the other hand, the sub-blocks are connected to each other through a common bit-line, shown as *SBL*.

Figure 4-8(b) shows the structure of the k^{th} column of a sub-block. It consists of 32, 8T bit-cells on a column. The WWL and RWL signals of those bit-cells are generated by their row circuits which are not shown in the figure. The local sensing is performed through the two VTC shifted inverters and a PMOS.

A typical read path of a 6T bit-cell based SRAM was given in Figure 1-11. As it was pointed out, read path is generally the critical path in SRAM memories [58]. This is due to the fact that the BL capacitance is a large capacitance since it connects to the drain nodes of all the PG transistors on a column and it is a long metal wire with a large wiring capacitance. So, discharging of the BL capacitance takes a long time and it affects the performance of the memory. Therefore, lowering the BL capacitance

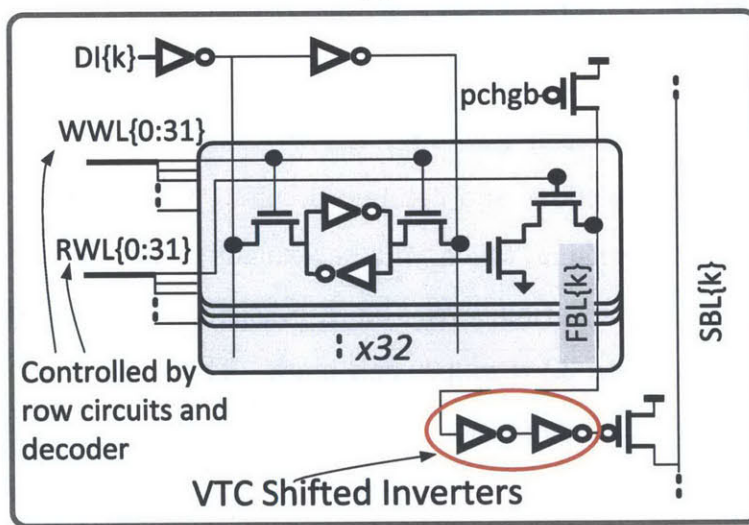
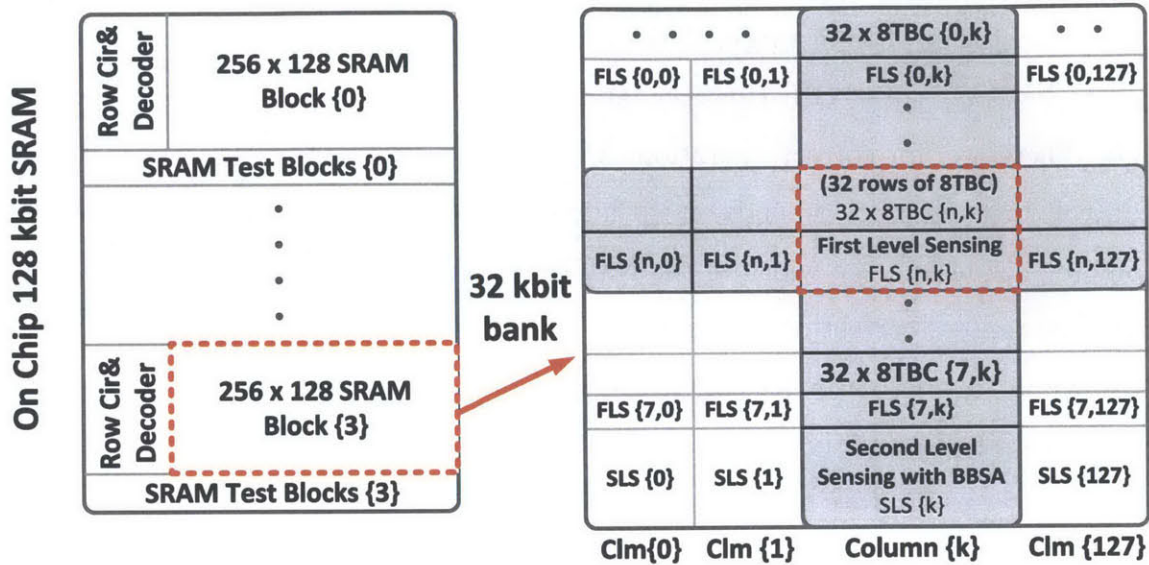


Figure 4-8: Organization of the 128 kb SRAM. (a) The two-stage sensing of the 32kb block. (b) The schematics of the k th sub-block.

would result into a faster operation. Furthermore, SA input offset compensation is an effective method to improve SRAM performance. This technique tightens the offset distribution of the SA and enables correct sensing with a smaller voltage differential appearing on the BL.

As pointed out earlier, since 8T bit-cell cannot be used with column interleaving, adjacent columns cannot share a SA. As a result, each column has its own SA which

makes the area of them to be more constrained. Therefore, it is harder to design a low-offset SA for an 8T bit-cell based memory compared to its 6T counterpart.

This performance degradation problem of the 8T bit-cell based memory is addressed with two design strategies in this thesis: (1) Using two-stage sensing scheme, and (2) Utilizing a new offset-compensated SA using body-biasing.

4.1.5 Two-stage Sensing

Conventionally, SRAMs employ a small signal sensing scheme where signal development on long bit-lines are amplified by complex sensing circuitry. In this thesis, a two stage sensing is used. The important modifications to the conventional sensing scheme are:

- The long bit-lines of the conventional SRAMs are replaced by: (1) Short first-level bit-line (FBL)'s which are connected to the bit-cells and, (2) Long second-level bit-line (SBL)'s that traverse the entire memory (Figure 4-8).
- The signal development on those bit-lines are sensed in two stages: (1) First-level sensing (FLS) which uses two inverters and a PMOS transistor, and (2) second-level sensing (SLS) which is performed by the proposed body-biased sense amplifier (BBSA).

In the two-stage sensing scheme, at the beginning of a read operation, the FBLs are precharged to V_{DD} . Precharging of FBLs is performed by the *pcghb* signal which is connected to the PMOS transistors as shown in Figure 4-8. Similarly, SBLs are discharged to zero using NMOS transistors with *pchg* signal which is not shown. The read operation starts with RWL signal assertion. Figure 4-9 shows the important signals of the single-sided read path during two back to back read operations.

1. The first clock cycle illustrates a read of data '0'. During this cycle, both FBL and SBL stay at their precharged values.
2. The second clock cycle illustrates a read of data '1'. During this cycle, FBL discharges to '0'. This is transferred to the SBL and it starts to charge up.

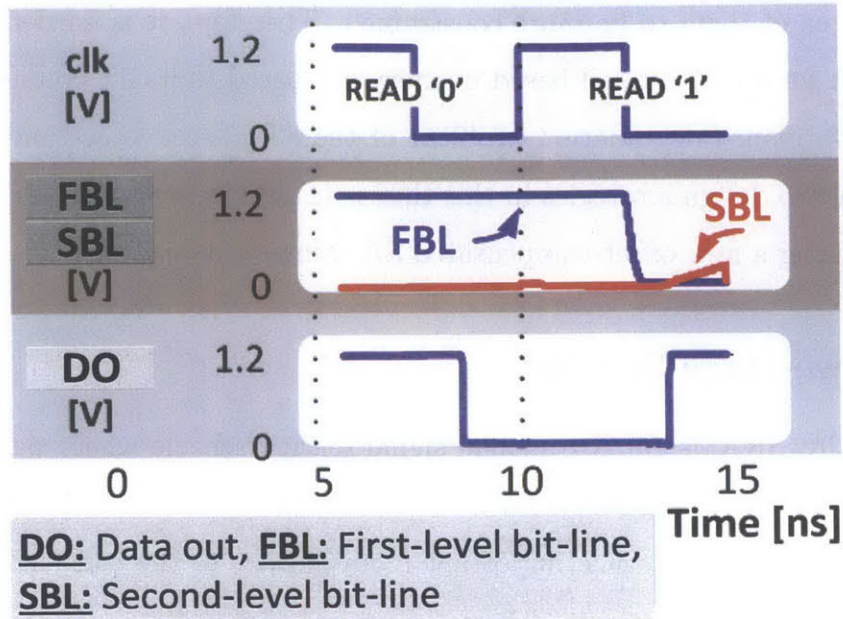


Figure 4-9: Two stage sensing signals during two back-to-back read operations.

In this work, FLS consists of two static inverters. It should be noted that using one inverter would be enough for sensing the signal development. However, having two inverters instead of one enables data-dependent read energy consumption. In this work, if the data that is being read is a '0', both FBL and SBL stay at their precharged values. On the other hand, if only one inverter is used, one of FBL or SBL would always need to discharge. Therefore, for applications where data stored in the SRAM statistically favors zeroes, this memory would result into up to 2× energy savings.

The FLS inverters are designed to have shifted voltage transfer characteristics where the inverter trip point, V_M , is shifted by a 100 mV in order to result into faster tripping. Figure 4-10 shows the effect of VTC shifting of the FLS inverters. The figure shows the ratio of read access time improvement (%) to the FLS area increase (%). This ratio should be maximized to have maximum read time improvement for minimum area overhead. For this design, 100 mV shift maximizes this ratio.

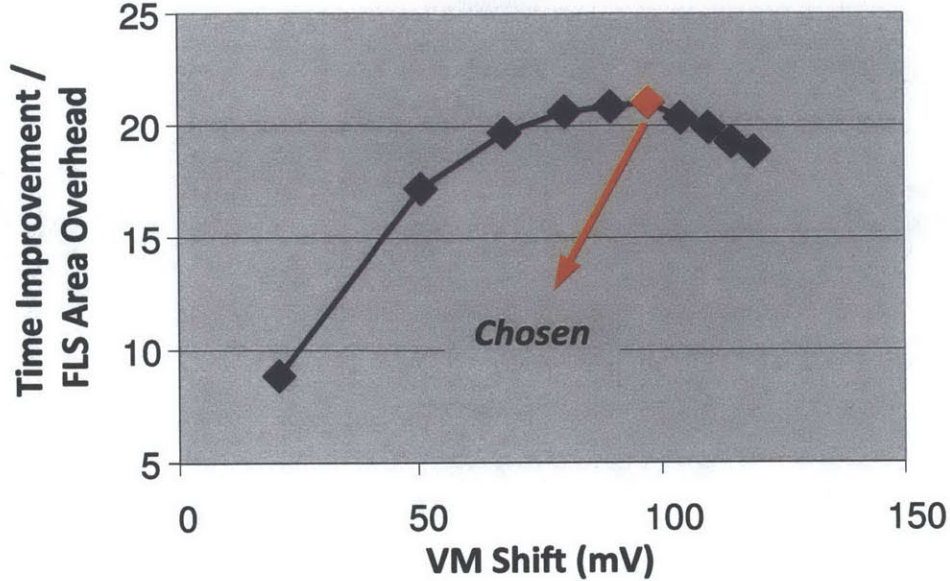


Figure 4-10: First level sensing (FLS) inverters use a 100 mV VTC shift to optimize time improvement vs. area overhead.

4.1.6 Sense-Amplifier Offset Reduction Using Body-Biasing

The ability of a SA to sense small signal input differences is limited by its input-referred offset. Offsets in SAs occur due to both global and local variations as explained in Chapter 1. The global variation effect can be mitigated using differential SAs. Since 8T bit-cell has only one read BL, it cannot work with the conventional differential sensing; so, large signal sensing is commonly used with 8T bit-cells. During large signal sensing, the BL needs to be discharged from V_{DD} to ground. Therefore, it tends to increase the energy overhead and access time. Due to this reason, in this design, a *pseudo-differential sensing* is used as SLS. SBL drives one of the inputs whereas the other input is driven by a reference voltage, REF.

One design approach that has been used to decrease the SA offset has been increasing the area of it. However, the fact that 8T bit-cells cannot share SAs stresses a general problem observed in deeply scaled technologies. Specifically, the size of the SA is not scaling due to the trade-off between their statistical offset and their layout area [83]. In this design, this trade-off is managed by using a body-biased SA design.

In Figure 4-11 (a), BBSA circuit implementation proposed in this work is shown. SA offset voltage can be either negative or positive depending on which side of the SA

or M2 is body-biased after calibration, since which one needs to be body-biased is not known in advance, the n-wells of M1, M2 and the rest of the PMOS devices have to be laid out far from each other. This makes the layout of BBSA challenging. By placing NMOS devices of the BBSA between different n-wells and carefully designing the layout of the circuit, total SRAM area overhead due to BBSA is kept below 3.5% compared to total SRAM area.

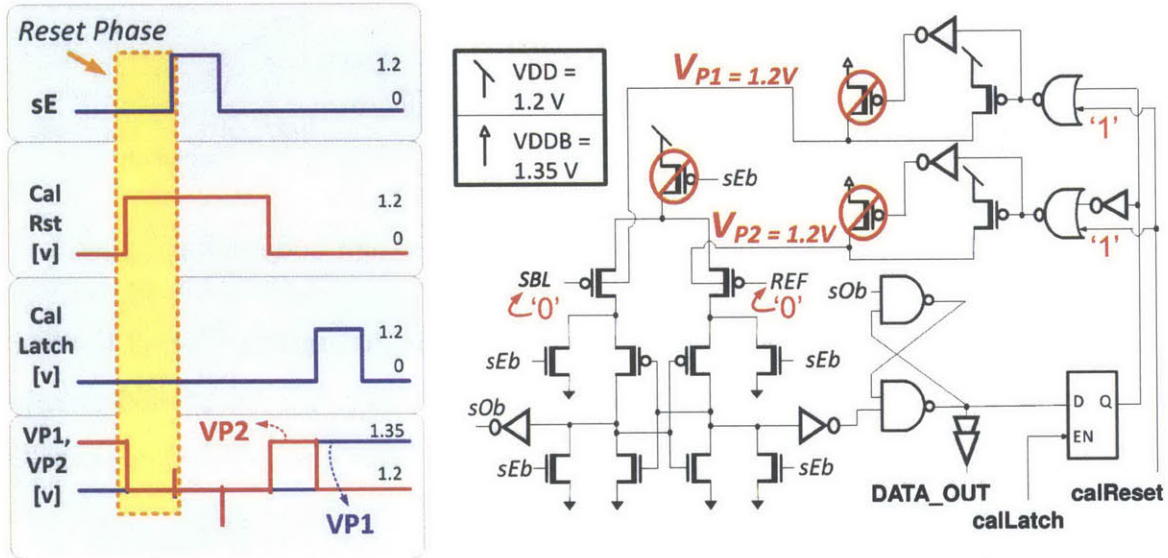


Figure 4-12: BBSA operation explained in detail - Reset phase.

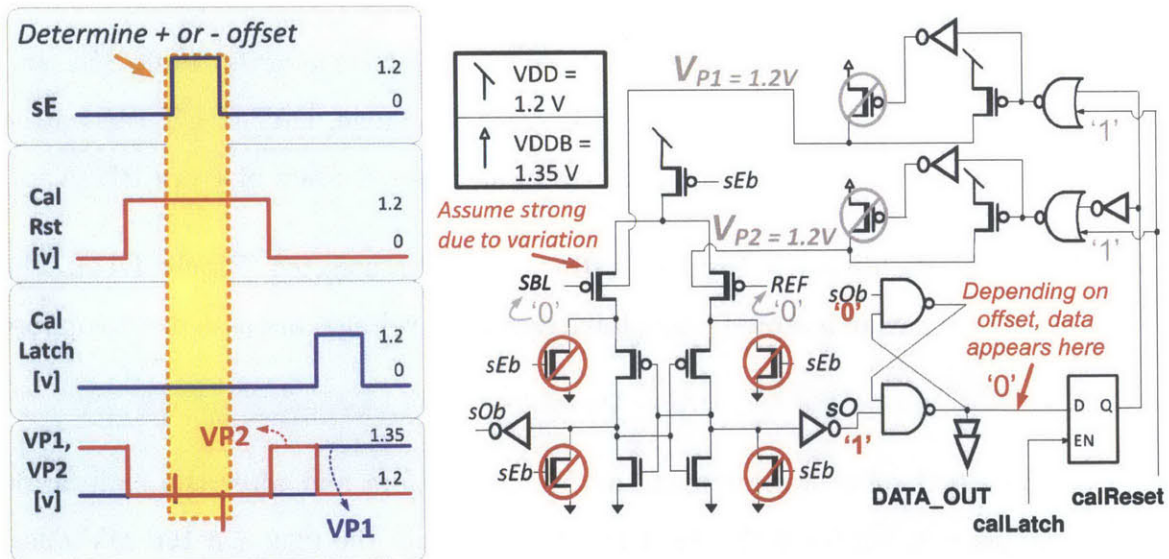


Figure 4-13: BBSA operation explained in detail - Determine offset phase.

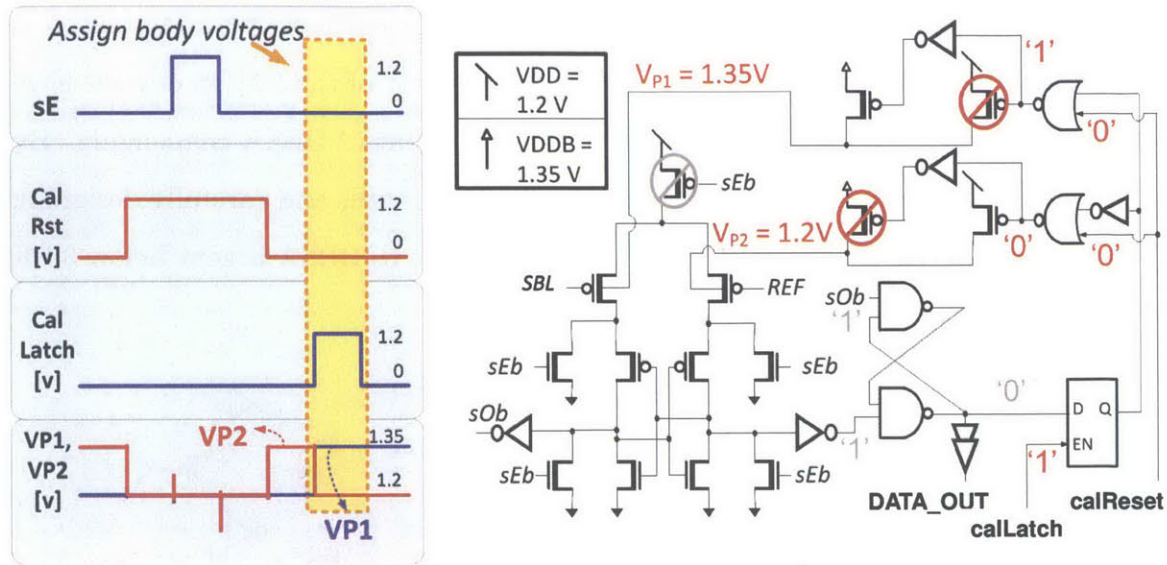


Figure 4-14: BBSA operation explained in detail - Assign body voltage phase.

The calibration process has to be performed once at the startup. For all the SAs in the design, it can be performed simultaneously. Therefore, this calibration process is simple (requires one latch and a few gates) and fast (<5 clock cycles). It can be summarized as follows:

1. During the *reset stage*, calRst is kept high to make VP1 and VP2 of all BBSA circuits equal to $V_{DD} = 1.2\text{ V}$ as shown in Figure 4-12.
2. During the *determine offset* phase, sense enable (sE) is asserted while SBL and REF are 0V as shown in Figure 4-13. After this stage, DATAOUT carries the information about positive or negative input referred offset of every BBSA.
3. Finally, at the *assign body* phase, calLatch signal is asserted (Figure 4-14). The offset information is stored into a latch and body voltages are assigned to oppose the offsets.

The measured offset variation of the 512 SAs before and after the calibration process is given in Figure 4-15. As it can be seen from the figure, a 190 mV offset span appears before calibration. A 150 mV body-biasing results into a 50 mV offset shift in the distribution per SA. The SAs with a negative offset are compensated in

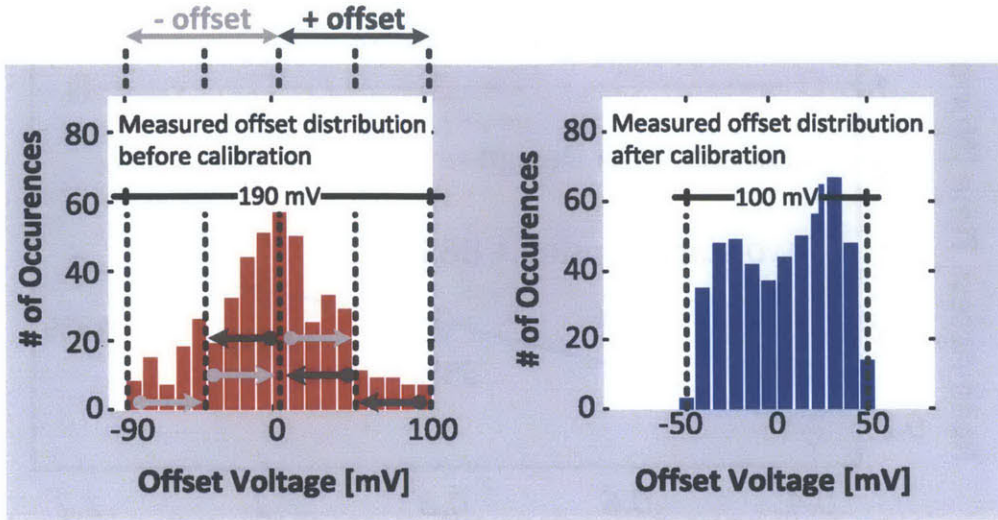


Figure 4-15: Measured BBSA input referred offset voltages before and after calibration.

Table 4.1: Comparison of the sense amplifier with offset compensation using body-biasing (BBSA) to recent work.

	This work	[60]	[62]
Meas. Offset Reduction	2×	5.75×	NA
Calibration Required?	Yes	Yes	Yes
Calibration On-chip?	Yes	No	Yes
Calibration Time	<5 clock	NA	NA
# of extra reference voltages	1	16	None
Calib. Circuit per SA	1 latch, a few gates	External	2 FFs, a few gates

the same fashion whereas the SAs with positive offset are compensated in the opposite fashion. For instance a SA with a -90 mV offset voltage ends up having a -40 mV offset after the calibration. Similarly, a SA with a 40 mV offset results into having a -10 mV offset. This way, the distribution is 2× tighter after calibration. Compared to an upsized SA that has a 100 mV input offset span, BBSA requires 1.5× smaller area.

Table 4.1 compares BBSA to recent work. BBSA calibration process is very simple since it requires only one latch and a few gates and it takes <5 clock cycles for

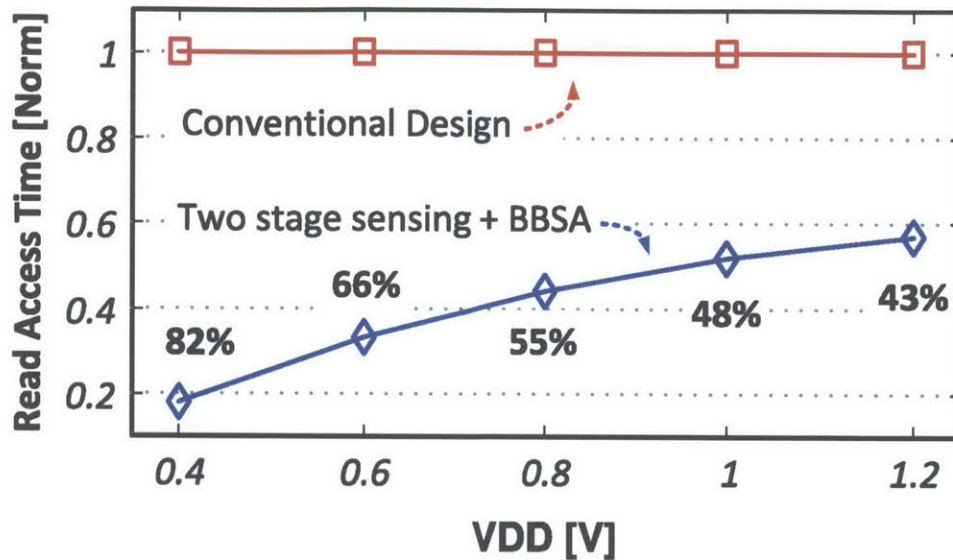


Figure 4-16: Improvement in performance by using two stage sensing and BBSA compared to the conventional design.

calibrating all SAs. Moreover, BBSA offset compensation can also apply to differential SAs since it requires controlling the body voltage rather than the REF input.

Figure 4-16 compares the statistical simulation results of read access time of the two-stage sensing and BBSA to conventional sensing. For the conventional case, a memory with 256 bit-cells per bit-line that uses small-signal strong-arm type SAs without offset compensation is assumed. The sensing scheme brings up to 82% read time improvement due to two-level sensing and offset reduction using BBSA compared to the conventional way.

4.1.7 Measurement Results

Ideas presented in this section are implemented in a 65nm low-power CMOS process. A die photograph of the test-chip is given in Figure 2-14 and Table 2.1 summarizes the features of this test-chip. Four SRAM macros are placed on a die with a total size of 128 kb on a die.

This work is explained in [72]. For testing purposes, the fabricated die are packaged into a -pin QFP ceramic package and a 4-layer test PCB board is designed.

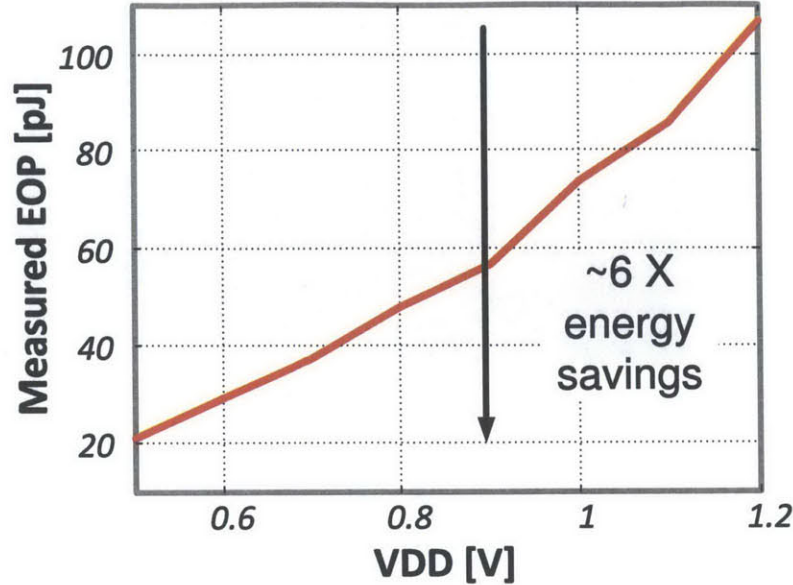


Figure 4-17: Measured energy per operation results vs. V_{DD} for the 128kb SRAM designed in 65nm CMOS.

Figure 4-17 shows the measured energy per operation numbers of the SRAM vs. V_{DD} . SRAMs operate from 1.2 V down to 0.37 V. It results into more than $5\times$ energy savings compared to operating the SRAM at nominal voltage of 1.2V.

4.2 A 6T Bit-cell Based SRAM in 28nm FD-SOI CMOS for Operation Below 0.5 V

Although 8T bit-cell is promising for low-voltage operation, it has many disadvantages compared to its 6T counterpart. Therefore, 6T design is still the conventional bit-cell for SRAM design. In this section, we will focus on methods to make the 6T bit-cell based SRAM arrays to work at low-voltages (down to 0.43 V) using peripheral assist circuits and scaled technology nodes like 28nm FD-SOI.

In this design, a high-density, $0.152\mu\text{m}$, 6T bit-cell is used. The SEM image of the bit-cell is shown in Figure 4-18.

The principle mechanisms used to enable low voltage operation of this SRAM are:

1. single-sided, cycle-by-cycle basis body-biasing for better write-ability



Figure 4-18: SEM picture of the $0.152\mu^2$ 6T bit-cell. This bit-cell is used in the 0.5Mb SRAM test-chip [7]. *Courtesy of STMicroelectronics*

2. short local BLs to minimize read disturbances

4.2.1 FD-SOI Technology Specifics

As transistor scaling reached 45nm and beyond, it has become harder to lower the supply voltage of the SRAMs due of increased transistor variation. To enable further scaling, new process technologies started to appear in SRAM design. At 45nm node, high-k technology became popular [70, 84]. At 28nm and beyond, FinFET [85, 86] and SOI transistors [87, 88, 89] started to show up in SRAM design.

FinFET (or trigate) a promising technology since it reduces short channel effect (SCE) and leakage currents, and enables further voltage scaling. However, it is a fundamental change over bulk CMOS since its transistors do not have a planar structure but has a 3D structure. Furthermore, in FinFET technology, transistor width is quantized and cannot be sized freely. This creates a major concern for bit-cell design for SRAMs since it gets harder to optimize the bit-cell for read and write operations.

An alternative technology is the Fully Depleted Silicon On Insulator (FD-SOI) which is a planar technology but can deliver reduced leakage and variation compared to bulk transistors. The cross-section of a typical FD-SOI transistor is shown in Figure 4-19. In FD-SOI, first, a thin layer of insulator, buried oxide (BOX), is positioned on top of the base silicon. Then, a very thin silicon creates the channel. The 28nm FD-SOI CMOS transistors are fabricated in a 7nm thin layer of silicon for

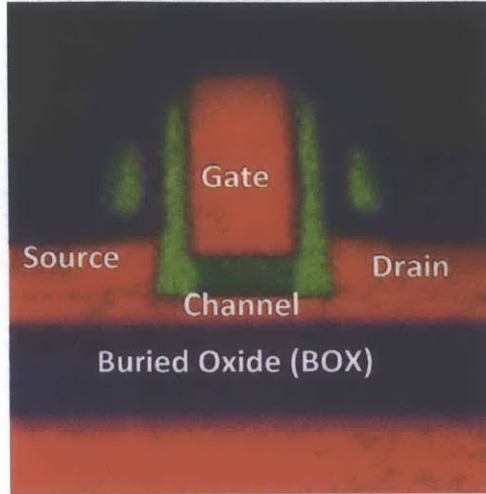


Figure 4-19: FD-SOI transistor cross-section [8].

the channel sitting over a 25nm BOX [90]. This is called the ultra-thin body and buried oxide (UTBB) FD-SOI process. Thanks to its thickness, there is no need to dope the channel, making the transistor fully-depleted [91]. FD-SOI provides a much better electrostatic control compared to bulk CMOS [92] since electrostatics are not controlled by channel doping [93].

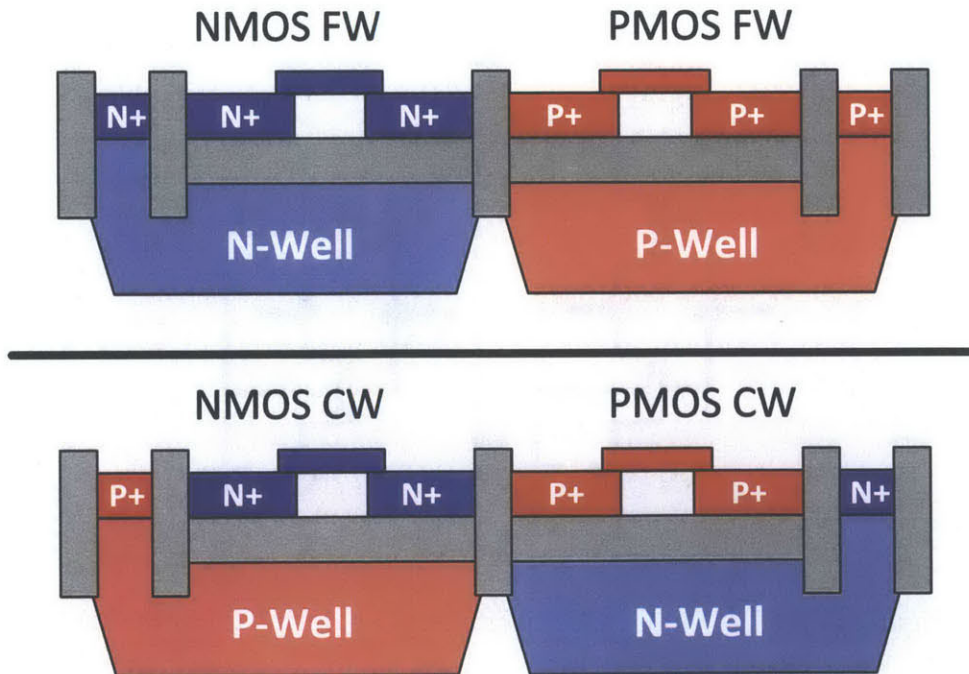


Figure 4-20: Flip-well (FW) vs. conventional-well (CW) transistor structures of FD-SOI.[7]

Moreover, FD-SOI technology enables control of the behavior of transistors not only through the gate but also through the *back gate*. The back gate control is similar to body-bias of bulk technology and achieved by polarizing the substrate. However, body-biasing is much more effective in FD-SOI [94]. A higher body factor of 85mV/V can be obtained compared to bulk technology which can achieve 25mV/V. Moreover, it is important to note that the body-bias range in bulk technology is limited to -300mV because of source-drain junction leakage and latch-up risk at higher voltages. However, UTBB FD-SOI technology enables an extended body-bias range from -3V reverse body-bias (RBB) up to 3V in forward body-bias (FBB). This provides a new knob in energy efficiency optimization, performance boosting, ultra-low functionality and leakage reduction [95].

In conventional well (CW) processes, such as bulk CMOS, PMOS transistors are manufactured on n-well; whereas, NMOS transistors are on p-well. The same CW structure can also be used in FD-SOI. However, in FD-SOI, a special flip well (FW) structure is also possible. By flipping the wells of CW, NMOS transistors are manufactured on n-well, and PMOS transistors are manufactured on p-well. The comparison of the transistor cross-sections which are fabricated in FW and CW processes are shown in Figure 4-20.

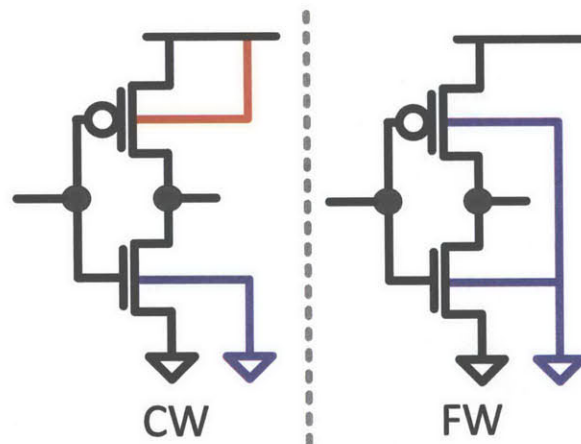


Figure 4-21: FW vs. CW body connections during no body-biasing.

The body connections of an inverter gate which is designed in FW and CW processes of UTBB FD-SOI are shown in Figure 4-21. The CW process is similar to the

bulk design where PMOS body is connected to V_{DD} whereas NMOS is connected to GND. On the other hand, in FW, both NMOS and PMOS body are connected to GND. In FW, in order to avoid the direct p-n junction biasing between the n-well and the p-well, the well bias of the PMOS needs to be connected to GND rather than V_{DD} . Additionally, the NMOS body voltage can be increased from GND up to 3V to decrease the NMOS V_T dynamically. This is particularly important for our design since selectively body-biasing the NMOS transistors is beneficial for write-ability. This will be explained in more detail in the following sections.

4.2.2 Low-Voltage Challenges of 6T Bit-cell Operation in FD-SOI

Compared to bulk technology, due to the elimination of channel dopants, LER and RDF variations in the channel are suppressed for UTBB FD-SOI technology. As a result of reduced variation, around 0.7V operation with a 6σ yield was projected for a 6T bit-cell [96, 97]. This is a significant improvement since SRAMs that are built using bulk technologies can barely achieve 0.85 V operation using a bulk process as can be seen in Figure 4-1. On the other hand, the minimum energy point of circuits generally lies in the subthreshold range [98]. Therefore, lowering the operating voltage of the SRAM circuits from 0.7 V down to below 0.5 V would not only result into better energy efficiency for the SRAM circuits but also a better system integration.

To analyze the low-voltage operation of the 6T bit-cell, the write-ability is analyzed under variation effect. bit-line write margin (WVBL) vs V_{DD} graph is shown in Figure 4-22. The simulation methodology is explained in Section 1.3. 10K MC analysis is used to evaluate the variation effect with the assumption of a Gaussian distribution to estimate μ and σ . As it can be seen from the figure, for a $6\mu/\sigma$ target, this bit-cell cannot work below around 0.75 V due to reduced write-ability.

Similarly, to quantify the bit-cell's robustness against read failures, SNM simulations are performed. In order to generate this graph, 10K MC analysis is performed. Figure 4-23 shows the read SNM μ/σ on y-axis and V_{DD} on x-axis. As it can be seen

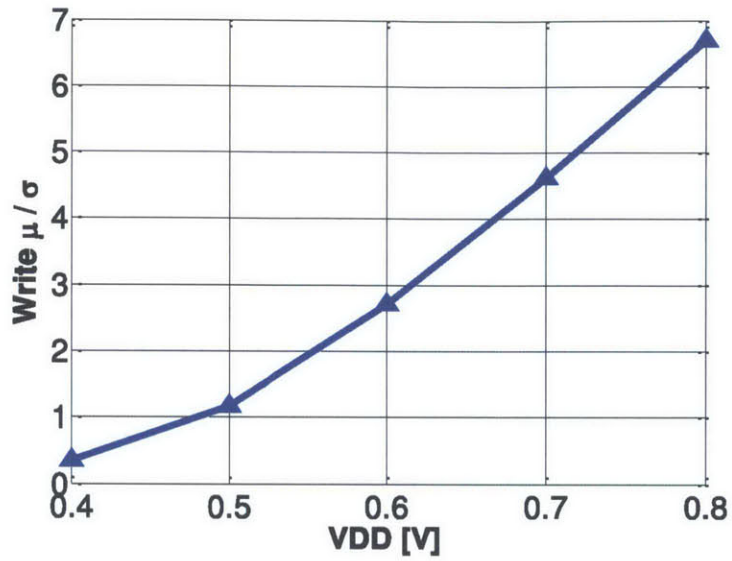


Figure 4-22: Write margin μ / σ vs. V_{DD} .

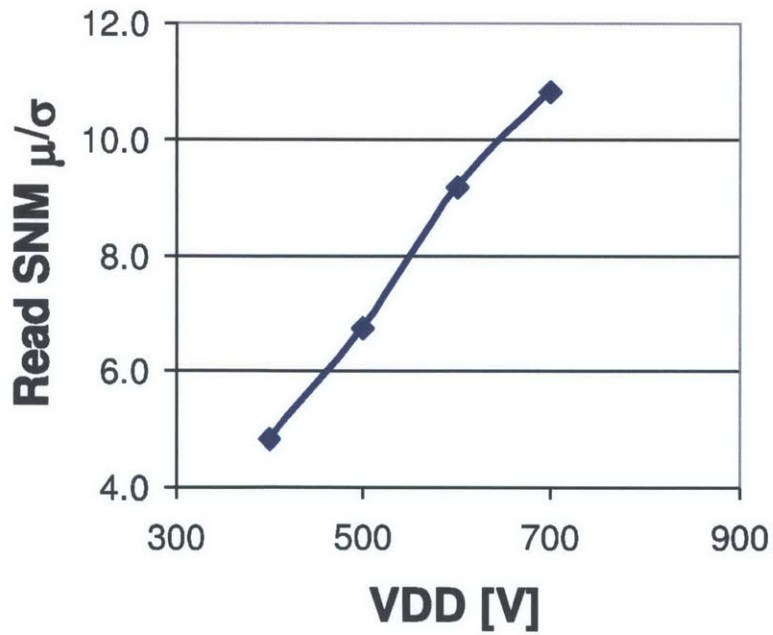


Figure 4-23: Read SNM μ / σ vs. V_{DD} .

from the figure, read operation is also expected to cause upsets under around 500 mV.

One challenge with the 6T bit-cell is that in general an assist technique that helps with the write operation make read-ability worse, or vice versa. For instance, increasing the WL voltage helps with write operation but it degrades read-ability. Therefore, the effects of the assist circuit to both operations need to be carefully considered.

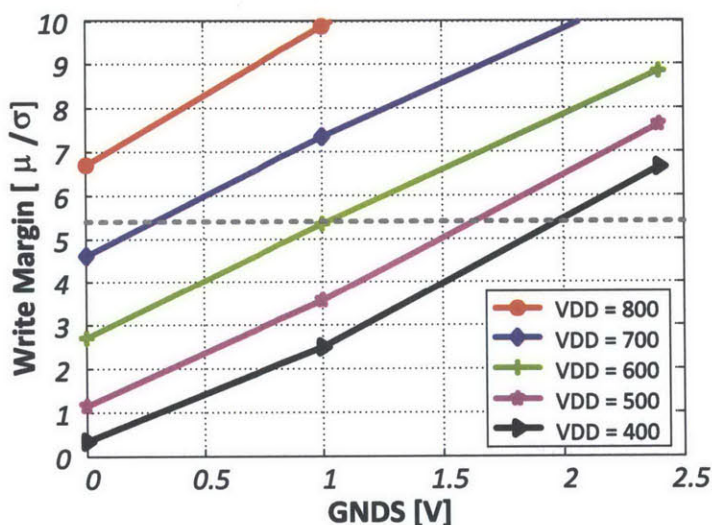


Figure 4-24: WVBL under body-biasing under worst case operating conditions.

4.2.3 Write-Assist using Body-Biasing

Bit-cell write operation was explained in Chapter 1 and the bit-cell under a write operation can be seen in Figure 1-12. Since for stable write operation, PG needs to be stronger than the PU transistor, decreasing the threshold voltage by body-biasing would increase the strength of the NMOS transistors and help write-ability.

Figure 4-24 shows the effect of increasing the body voltage of the PG transistors during a write operation to write-margin. Since increasing the n-well of the NMOS transistors decrease their V_T , this increase their write-ability. The figure shows the WVBL margin μ/σ on y-axis and the body voltage of the NMOS transistors (GNDS) on x-axis. The same analysis is performed for 5 different supply voltages (From 800mV down to 400mV with 100mV steps). For every point on the curves, 10K MC analysis is performed. As it can be seen from the figure, body-biasing the PG

transistors helps write-ability.

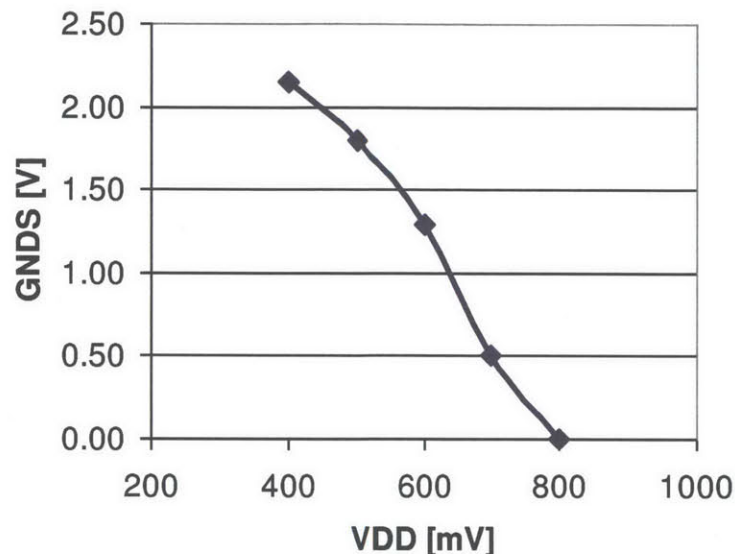


Figure 4-25: body-bias required for voltage scaling.

The required body voltage for a 6σ write-margin is given in Figure 4-25. This simulation assumes the worst case corner and temperature. As it can be seen from the figure, in order to make the bit-cell work down to 0.4 V, over 2V body-biasing would be required under the worst case assumption. This much body-biasing would not be possible with a bulk technology, but it is possible with FD-SOI.

Although body-biasing the PG transistors help write-ability, it degrades read stability. Figure 4-26 shows the effect of body-biasing to read SNM. As it can be seen from the figure, read-upsets are expected to increase due to body-biasing. Based on Figure 4-25, a body voltage larger than 1.5 V is required for write-ability constraint. However, at this body voltage and V_{DD} , read-upsets start to become the bottleneck.

In order to overcome the problem of degraded read-stability due to body-biasing, the technique that is proposed in this thesis is to perform body-biasing only during a write operation. This way, during a read operation, the PG transistors are not body-biased and read-ability is not degraded. However, this does not solve the problem of the *half-selected* bit-cells during a write operation. A typical 6T bit-cell based memory with a 2 to 1 column multiplexing was shown in Figure 1-11. During a write

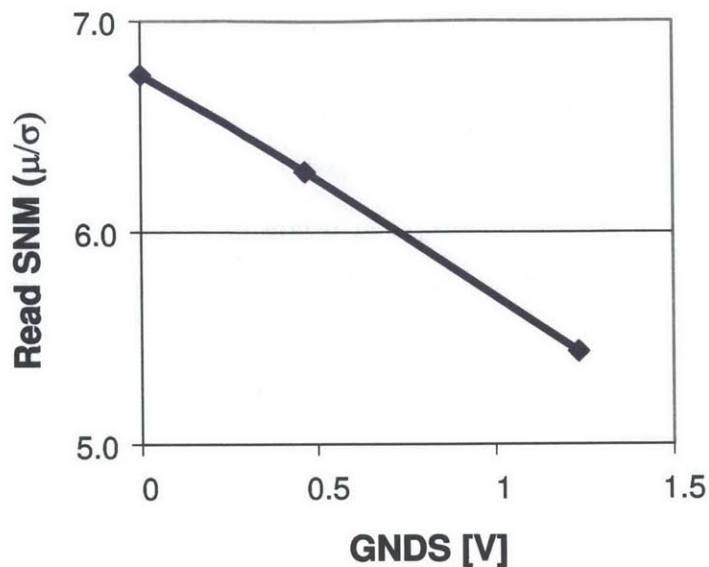


Figure 4-26: Read SNM under body-bias ($V_{DD}=500\text{mV}$).

operation of an SRAM with column multiplexing, the bit-cells that are on the active row, but are non-selected, are *half-selected* and they are under read-stress. Therefore, body-biasing their PG transistors during this period might result into read-failures.

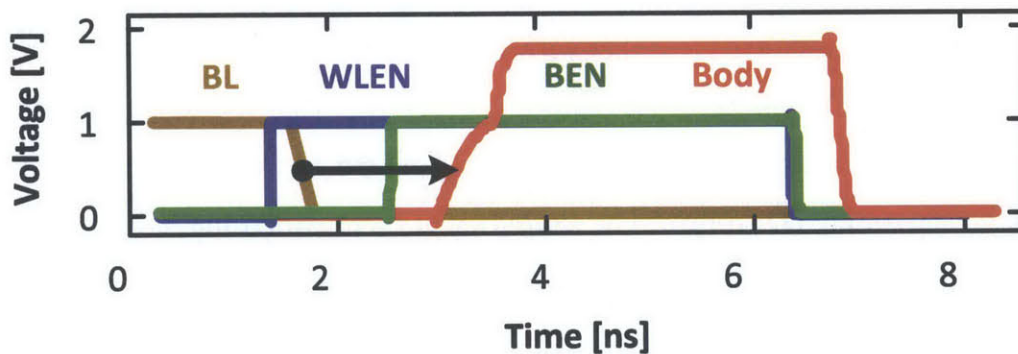


Figure 4-27: body-biasing starts after BL discharges.

To overcome this problem, the body voltage can be biased after enabling the WL. This way, the BL voltages would already discharge and read-upset problem would not occur for the half-selected bit-cells. This is performed in our test-chip by using different signals to enable body-biasing and WL (BEN and WLEN respectively). A typical operation is shown in Figure 4-27. As it can be seen, BEN is asserted after

WLEN to make sure that BL's are discharged before body-biasing starts.

Although, the solution to the read-upset problem is relatively simple to overcome, another obvious drawback of cycle-by-cycle body-biasing is the fact that it requires charging up the body capacitor during every write operation. Therefore, for net energy savings, this energy consumption needs to be minimized. Three techniques are used to decrease the energy consumption due to body-biasing:

1. Single-Sided Body-Biasing
2. Stepwise Charging of the Body Capacitor
3. Using Shorter BLs

In the rest of this section, those three techniques will be explained in detail.

Using a Single-Ended Write Operation

A 6T bit-cell that is designed using a FW process can be seen in Figure 4-28. Thanks to the FW technology, the two NMOS transistors can be connected to separate n-wells. It should be noted that in the original SRAM design, the body voltages of all the NMOS transistors are shorted together with the help of a simple metal connection in the edge-cell. The only modification required is removing that metal connection. This way, the body voltages (shown as GNDSA and GNDSB in the figure) of all the bit-cells can be separately controlled. The simplicity of this method made it possible for us to be able to perform this modification without changing any dimensions of the bit-cell or the edge cells.

Since the write operation starts from the side where a data zero is being written, body-biasing the PG at that side only, results into almost the same write-ability improvement compared to body-biasing both PGs. For instance from the figure, if during the write operation, BL is driven to '0' and BLB is '1', body-biasing GNDSA only would have almost the same effect compared to body-biasing both GNDSA and GNDSB. Based on the simulations, body-biasing both sides improves the μ/σ by only 0.1 compared to single-sided body-biasing. However, single-sided body-biasing

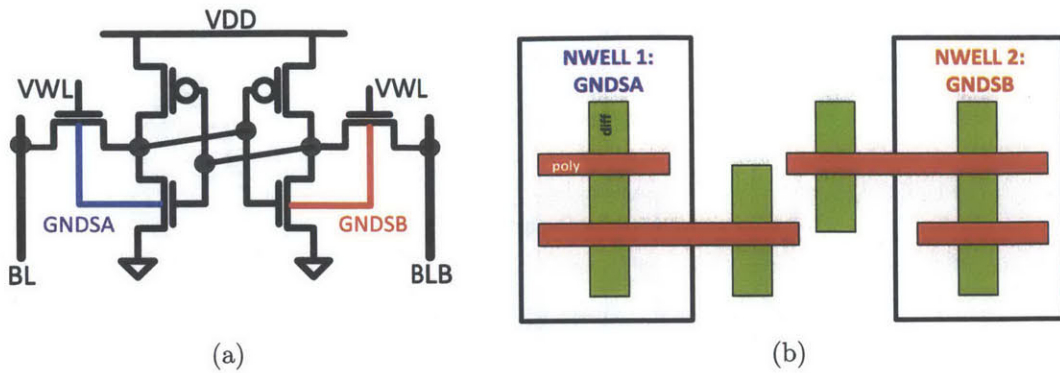


Figure 4-28: (a) Schematics, and (b) layout of the 6T bit-cell in the FW FD-SOI process. Layout is not drawn to scale.

requires half the energy compared to body-biasing both sides; with the trade-off of increased complexity.

Stepwise Charging of the Body-Capacitor

The second method proposed in this thesis for reducing the energy dissipation due to body-biasing is using a stepwise charging [99], which is an inductor-free form of adiabatic charging [100]. This method is illustrated in Figure 4-29. In this figure, a capacitance is being charged from a bank of voltage supplies with uniformly distributed voltages. To charge the capacitor, the supplies are switched on in ascending order, until the load reaches the final voltage. Here the voltages are given by $V_i = i \cdot V/N$. Stepwise charging results into reducing the dissipation by N .

In 28nm FD-SOI technology, the nominal supply voltage is 1V. To achieve operation below 0.5V, up to 2V body-biasing can be required as can be seen from Figure 4-25. Since the 1V supply is already required for the I/O, this intermediate voltage is used for stepwise charging. This way, the body voltage is not directly charged up to the high voltage but it is charged up in two steps: 0V to 1V; and 1V to VHIGH. This reduces the energy dissipation by around $2\times$.

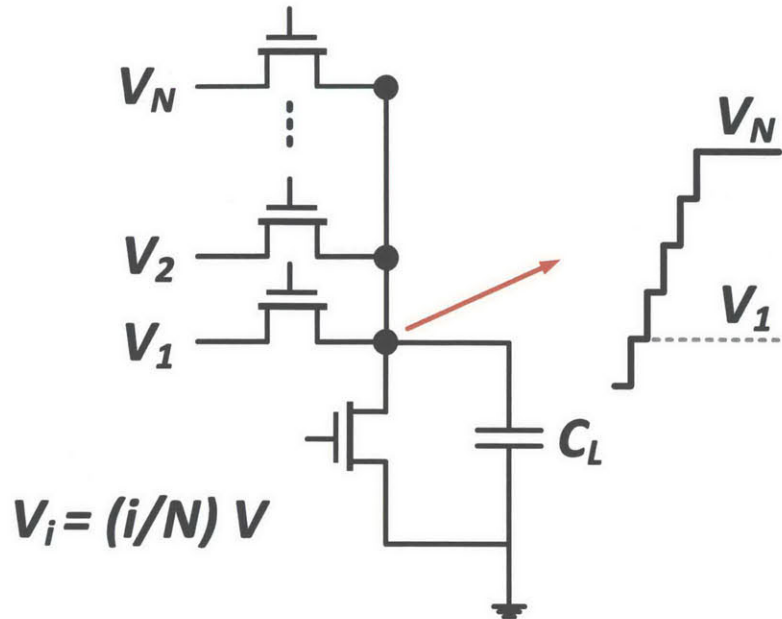


Figure 4-29: Stepwise adiabatic charging decreases energy dissipation by N .

Shorter local BLs

The third technique proposed to reduce the energy dissipation due to body-biasing is using a hierarchical BL scheme. The memory organization used in the test-chip prototype is given in Figure 4-30. The 0.5Mb memory is designed as 4 blocks where each block is 1024×128 bits. Within one block, a total of 32 sub-blocks exist which are designed in a 32×128 memory organization. After 8 sub-blocks, there is a *buffers* block to buffer control signals that are routed in a column-wise way. Each sub-block has a dedicated *local sense* block that handles the first step of data sensing. In addition to the local sense, there is a *body control* block per sub-block. Every 32×128 local array has a total of 129 GNDS body connections. All the GNDS voltages are controlled by one body control block.

There are two reasons for using the hierarchical BL scheme:

1. To have smaller GNDS capacitances to reduce energy consumption due to body-biasing.
2. To have smaller local BL capacitances to reduce read disturb.

Additionally, during body-biasing, since the V_T of the NMOS transistors are re-

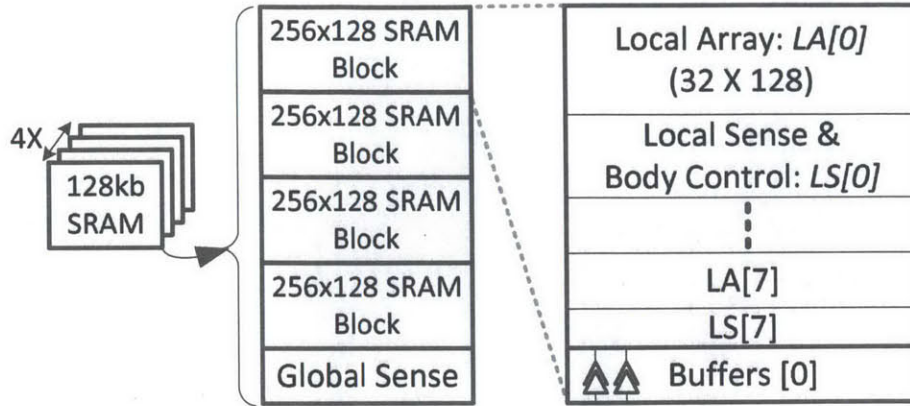


Figure 4-30: The memory organization of the 0.5Mb SRAM design.

duced, the leakage per sub-block increases. Dividing the memory into a total of 128 sub-blocks results into approximately $3.9\times$ smaller leakage compared to body-biasing 128 blocks together assuming body-biasing to 1V.

4.2.4 Memory Building Blocks

In this section, the design of the important memory blocks will be explained. The blocks that will be investigated here are the local sensing circuit, body-voltage control circuit, replica circuit, timing circuit, and built-in self test (BIST).

Local Sensing Circuit

The local sensing circuit schematics is shown in Figure 4-31. In this design, a column multiplexing ratio of 2:1 is used and handled by the column select signal (cSel). This local sensing circuit handles both read and writes. Both the local and global sensing uses large signal sensing. At the end of a read or write operation; BL, BLB, m, and mb are precharged to V_{DD} through their precharge transistors. The global bit-line (GBL) signal is also precharged to V_{DD} through a precharge transistor which is not shown in the figure.

At the beginning of a read operation, precharge signal (pchg) turns OFF; and, BL, BLB, m and mb nodes start floating. When WL is enabled, depending on the stored data inside the bit-cell, either BL or BLB is discharged. The selected column's

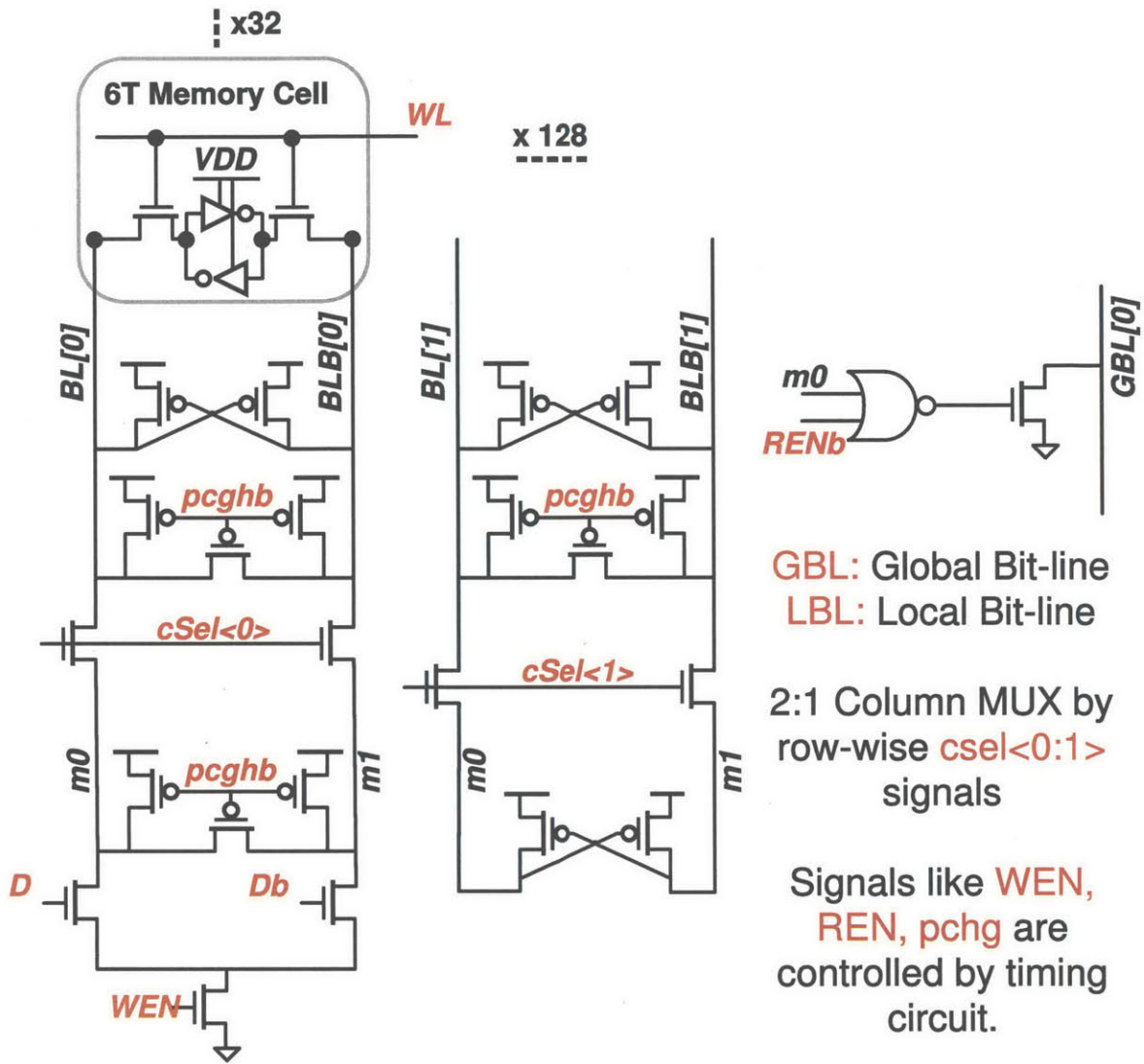


Figure 4-31: The schematics of the local sense circuit.

BL and BLB voltages are passed to the m and mb nodes. Then, REN is asserted which starts the global sensing through GBL. Using the two inverters in the global sensing circuit (which is not shown), the data is sensed.

The important read signals are shown in Figure 4-32. This simulation is performed for two back to back read operations. During the first clock cycle, a data '1' is read. On the next clock cycle, a data '0' is read. As can be seen in the figure, the signal development on the BL and BLB are passed to GBL. After the data appears at the output (DO), BLs and the GBL are precharged to their initial values ('1' and '0', respectively).

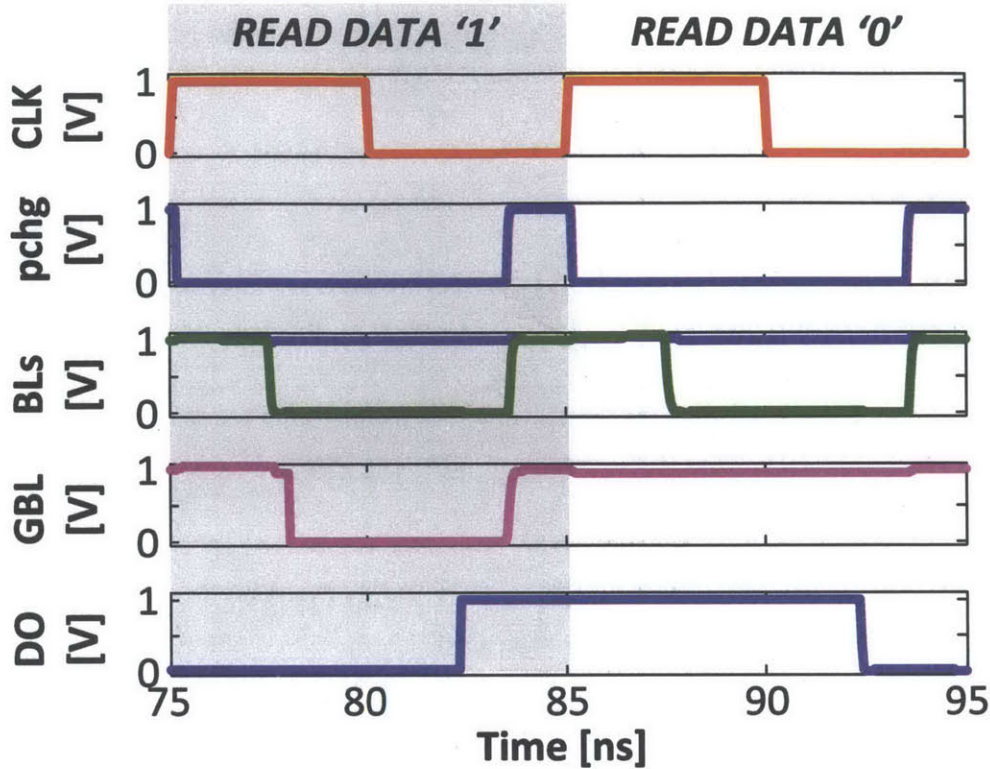


Figure 4-32: Important signals of the SRAM during two read operations.

At the beginning of a write operation, pchg signal turns OFF; and, BL, BLB, m and mb nodes start floating. Then, depending on the data to be written, data (D) and inverse of the data (Db) are asserted. When WEN is asserted, the data is passed to the selected column by the cSel. At this time, the neighboring column is half-selected and is under pseudo-read operation. Afterwards, the WL is asserted and the data is written into the selected bit-cell.

Body Control Circuit

Figure 4-33 shows the body control circuit. For every GNDS body connection, one body control circuit is used. In this design, during a write operation the body voltage can be:

1. charged up to V_{HIGH} through stepwise charging,
2. charged up to V_{DD} directly,

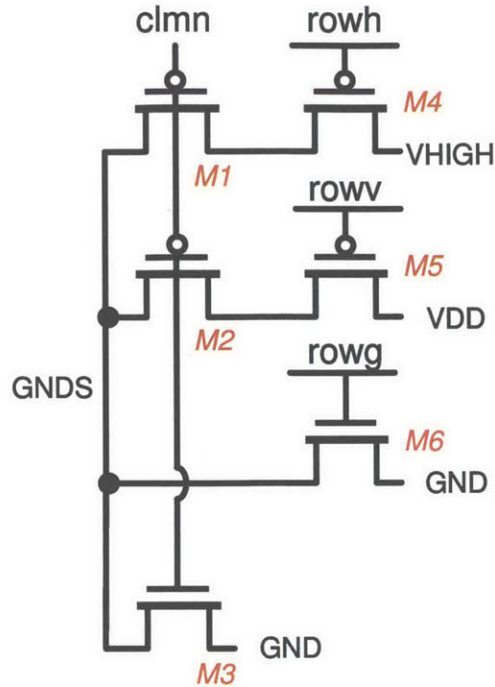


Figure 4-33: The schematics of the body control circuit.

3. can stay at GND if no body-biasing is selected.

The body voltage is controlled through four control signals: *clmn*, *rowh*, *rowv* and *rowg*:

- The *clmn* signal is routed column-wise. As mentioned, the single-sided body-biasing happens depending on which side of the bit-cell data zero will be written. Depending on that, the GNDS either needs to be charged up or it should stay low. (1) If it needs to stay low due to data being one, *clmn* signal is '1'. This way, M3 turns ON and M1 and M2 are OFF. (2) If it needs to be charged up, *clmn* signal is '0'. This way, M3 is OFF and M1 and M2 are ON. Depending on the row-wise signals, *rowg*, *rowv*, and *rowh* GNDS is charged up or stays at GND.
- The *rowg* is enabled (= '1') during write operation if no body-biasing is desired.
- The *rowv* is enabled (= '0') to charge the GNDS to 1V.
- The *rowh* is enabled (= '0') to charge the GNDS from 1V to V_{HIGH} .

Replica Circuit for Self-Timed Body Control

For the stepwise charging of the body capacitor to V_{HIGH} , first, $rowv$ needs to be enabled. Afterwards, when the body voltage, GNDS, reaches $\approx 1V$, $rowv$ needs to be disabled and $rowh$ needs to be enabled. The timing of this second stage is ambiguous and hard to predict. Therefore, in this design a self-timing replica circuit is used that generates a ready signal when GNDS voltage reaches a value close to 0.8 V.

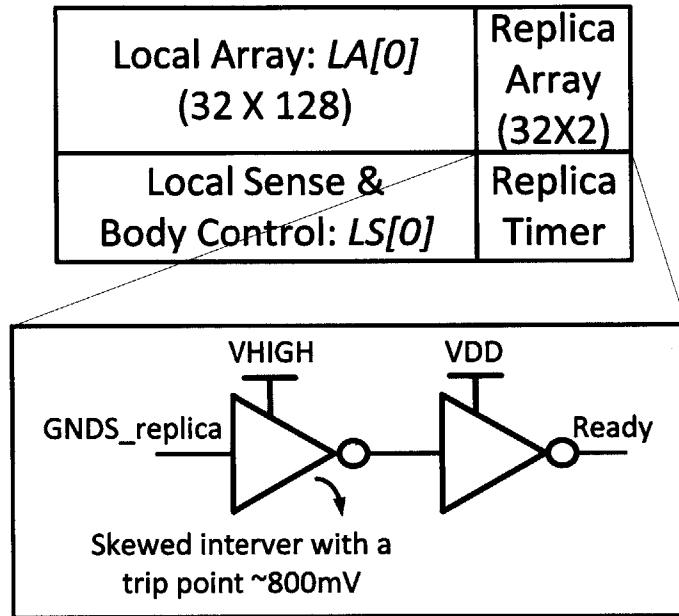


Figure 4-34: The schematics of the replica circuit.

Figure 4-34 shows the schematics of the replica circuit and its timer. Next to the local array, a 2b wide replica array is placed. Since there are two bit-cells in a row, there are three body connections, GNDS(0:2). The one in the middle, GNDS(1), or in other words, GNDS_replica, is used for creating the timing signal, *ready*. The *clmn* signal corresponding to GNDS_replica is always kept '0' so that during a write operation, independent of data, GNDS_replica will always start to charge up to V_{DD} . When GNDS_replica reaches a value close to 0.8V, the output of the first inverter in the replica timer circuit toggles to '0' and *ready* signal toggles to '1'. When *ready* turns ON, $rowv$ is disabled and $rowh$ is enabled. Then, GNDS starts charging up from V_{DD} to V_{HIGH} .

The effect of mismatch on the replica circuit would vary the timing of the *ready*

signal. Although this can have an effect on the total savings of this technique, since the replica circuit is only needed per sub-block, variation effects are not measured to be significant. Since we are using only two inverters to generate the *ready* signal, the replica timer circuit is simple such that it requires a small area overhead. This overhead is around 1% compared to the sub-block.

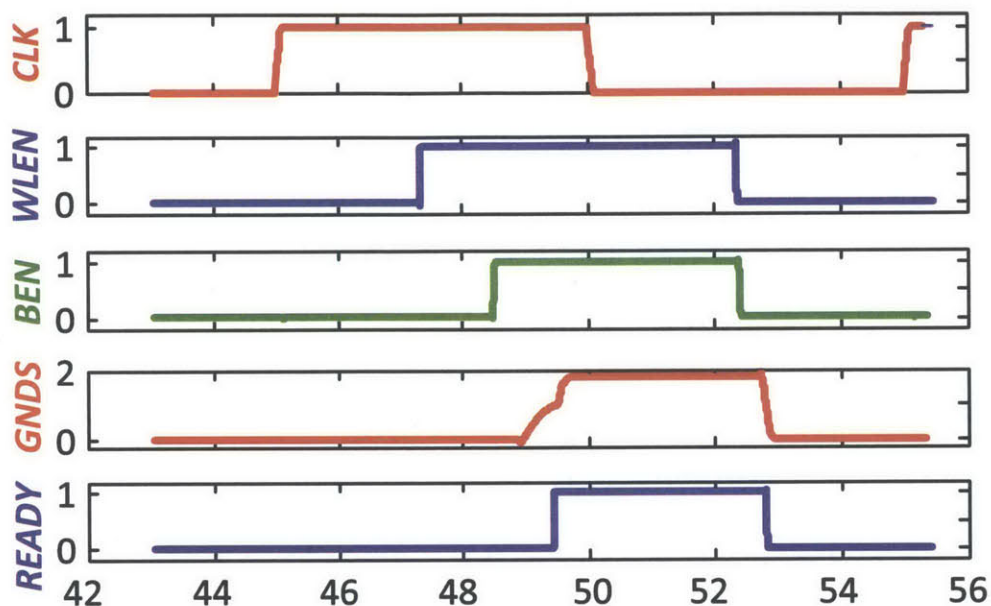


Figure 4-35: The signals used for stepwise body charging.

Figure 4-35 shows the important signals that are used for stepwise body charging during a write operation. The process works as the following:

- The *CLK* signal represents the clock to the SRAM. After the positive edge of the clock, word-line enable (*WLEN*) signal turns on, signaling the WL of the selected row to turn on. The delay between *CLK* and *WLEN* has to be larger than the worst case decoder delay.
- After some delay, body enable (*BEN*) signal turns on. During this delay period, the local BLs of the half-selected bit-cells discharge. This way, they are robust towards read-upsets. With *BEN* signal assertion, *GNDS* starts to charge up to 1V.
- When *GNDS* reaches around 800 mV, *READY* signal is asserted by the replica circuit. This results into *rowv* to be disabled and *rowh* to be enabled.

- For the rest of the cycle, the body voltage stays at its high value.

BIST Circuit

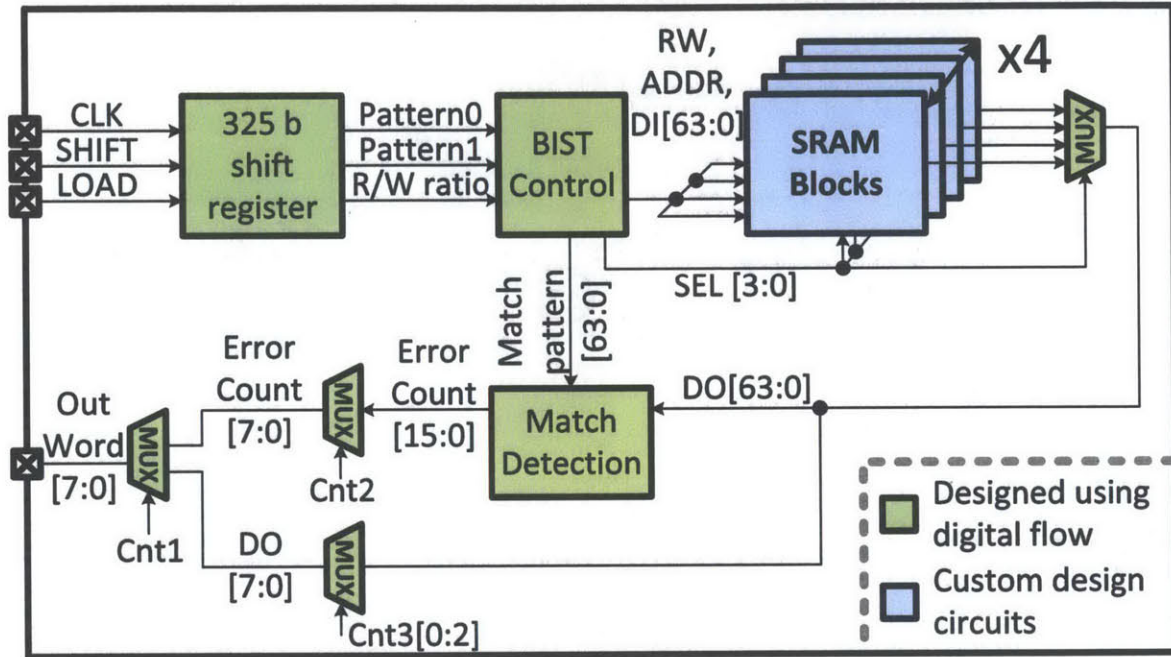


Figure 4-36: Block diagram of the BIST circuit for self-testing.

For testing purposes, a built-in self test (BIST) circuit is designed using the digital design flow. The block diagram of the BIST circuit is shown in Figure 4-36. The BIST circuit requires only 3 inputs and 8 outputs to generate different testing structures for the custom block SRAMs. As the first step, a 325b shift register is loaded in 325 clock cycles. Then, depending on the configuration of the shift register, the inputs such as read/write input (RW), address input (ADDR) or input data (DI) are generated within the chip. The test circuit can automatically count the number of bit-errors that are accumulated or it can directly pass the output of the SRAMs. Having the BIST circuit makes the testing simpler and it reduces the number of required IO pads in the pad ring.

4.2.5 Improving Read Stability Using Hierarchical Bit-lines

Conventional SNM stability metric cannot capture the dynamic behavior of the bit-cell during a read operation [4, 63, 57, 64]. The WL, BL, and internal storage node behavior needs to be taken into account for a more realistic read margin estimation. Several works have recently investigated how the dynamic operation of the bit-cell would affect bit-cell stability.

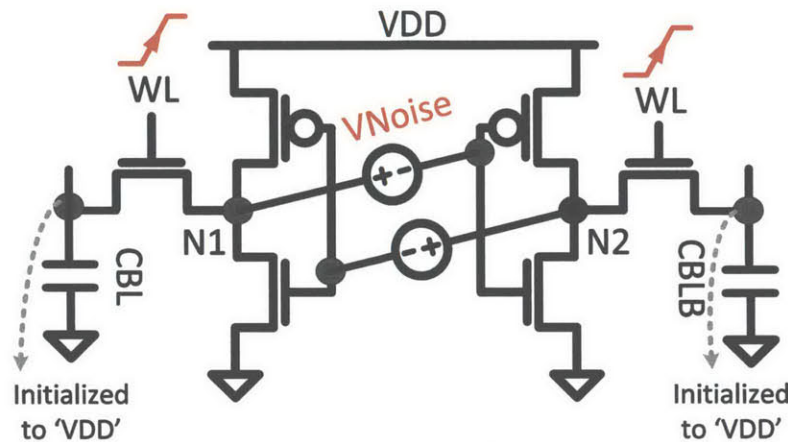


Figure 4-37: The simulation setup for the dynamic read stability margin calculation of the 6T bit-cell.

During the SNM simulation, a DC analysis is performed. It assumes the WL pulse width is infinite and that the BL and BLB are actively held at V_{DD} . Therefore, it gives pessimistic results. In reality:

- WL pulse is not infinite and the read disturb might not have enough time to flip the data stored in the bit-cell.
- BLs are not kept high during a read operation but they are floating. Therefore, the BL capacitance has a strong impact on the read stability.

Conventional SNM can overestimate the probability of read flip failure by 6 orders of magnitude [101]. Therefore, several dynamic read margin definitions are proposed to overcome this issue [102, 103, 56, 104].

Figure 4-37 shows the dynamic read stability simulation setup based on transient simulation. The margin is defined as the noise voltage that results into the storage

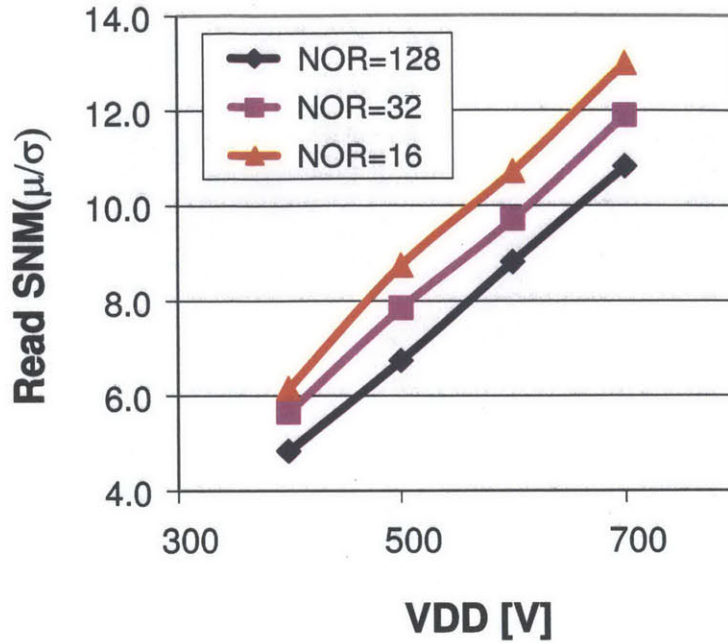


Figure 4-38: Dynamic read margin simulation results with NOR=8,32 and 128.

node to flip during a read operation. Here the two DC sources that are in parallel are the noise sources (V_{Noise}). The CBL and CBLB are the bit-line capacitances that are generated using layout extraction. The V_{Noise} is parametrically swept by 10 mV step sizes. A coarse-to-fine three-step-search is also possible for a better estimation [105].

The dynamic read simulation results for the 28nm FD-SOI process is shown in Figure 4-38. Here, the read stability numbers under number of rows (NOR) = 8, 32 and 128 are showed. As it can be seen, read stability increases as the number of rows is reduced. This is due to the fact that the bit-line capacitances are smaller for smaller NOR which results into a shorter time of read-disturbance. In this work, 32 bit-cells per bit-line is chosen to ensure read-robustness below 0.5V.

4.2.6 Test-chip Results

Based on the ideas explained, a 0.5Mb 6T bit-cell based SRAM test-chip prototype is fabricated using a 28nm FD-SOI process. Figure 4-39 shows the die photo of the chip. The die size is $1.8 \times 1.8mm^2$.

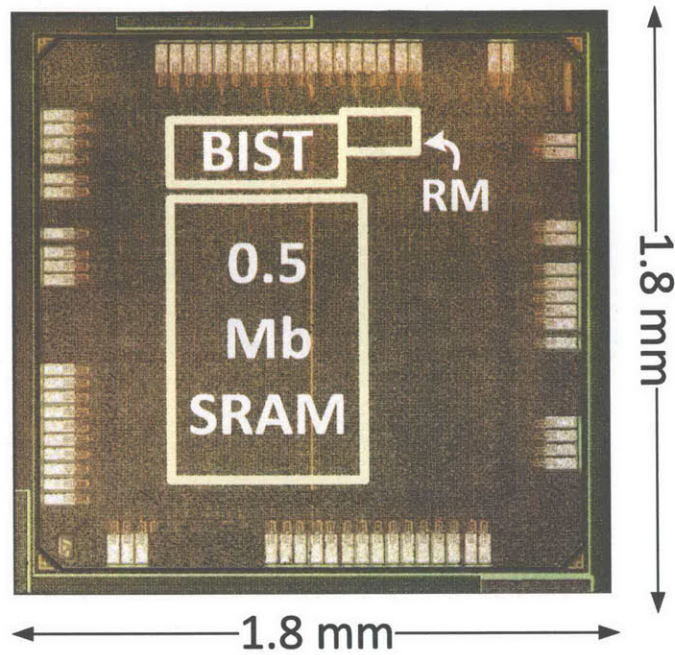


Figure 4-39: Die photo for the 28nm FD-SOI SRAM.

Table 4.2: Test-chip Specifications.

Technology	28nm FD-SOI
Bit-Cell Size	6T FW ($0.152\mu\text{m}^2$)
Chip Area	$1.8 \times 1.8 \text{ mm}^2$
Number of Pads	72
IO & Nominal Voltage	1.8 & 1.0 V
SRAM Capacity	0.5 Mb
SRAM Block Organization	4×32 sub-blocks
SRAM Sub-block Organization	$32 \text{ words} \times 128\text{bits/words}$
V_{DD} Range (measured)	1.0 - 0.43 V
Bit-cell Topology	6T
Column-Select Ratio	2:1
Body-Biasing Configurations	0V, 1V, 1.8V

Table 4.2 shows the specifications of the test-chip. The total size of the memory is 0.5Mb. There are four 128kb SRAM blocks and each block has 32 sub-blocks. The sub-blocks have 128 columns and 32 rows per column. A 2 to 1 column multiplexing ratio is selected for improving array efficiency and robustness against soft-errors.

The test-chip is implemented using a 72-pin wire bonding pad ring. This pad ring has only 14 I/O pads and the rest are dedicated for supply voltages. The small number of I/O pads necessitates a careful testing strategy. In this design, a BIST circuit is used that is built-into the test chip. This circuit enables automated testing and only requires 3 inputs and 8 outputs. Its gates are scattered around the SRAMs and it consists of around 2.5K gates.

To improve read-ability, 32 bits per LBL are used. Local sensing circuit accounts for around 14% of the sub-block area. However, this way, the global sensing can be done using only two inverters, therefore the effective area overhead is smaller. To improve write-ability, cycle by cycle body-biasing is used. The GNDS control transistors take up around 10.4% area. In addition, the replica circuit is used to create the step-wise charging of the body voltage and it requires around 1% of the sub-block area. The SRAMs are measured to be operational down to 0.43V.

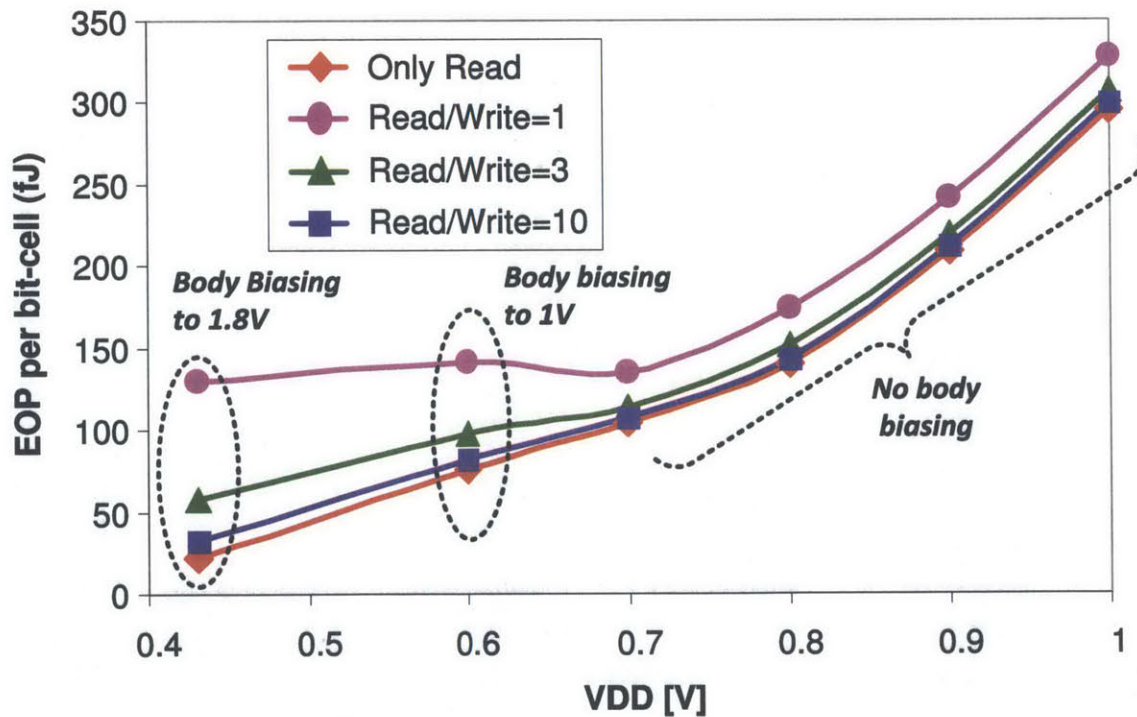


Figure 4-40: Measured energy consumption vs V_{DD} with and without write-assist for the 0.5Mb SRAM test-chip designed in 28nm FD-SOI technology.

Figure 4-40 shows the energy vs. V_{DD} graph of the memory for various read and

write ratios. This memory achieves operation down to 0.63V without body-biasing. To achieve operation down to 0.43V, it requires a 1.8V body-biasing. Since body-biasing is only performed during a write operation, the energy savings depends on the read to write ratio. This ratio is expected to be higher than 2 for many realistic applications. The measurement results can be summarized as follows:

- If the number of read and write operations are equal, the energy savings achieved by voltage scaling is not enough to compensate the energy consumption due to body-biasing. Thus, operating the memory at 0.43V results into a 14% higher energy consumption compared to operating it at 0.63V.
- If for every write, there are two read operations, we would achieve 37% energy savings.
- If this ratio is 3, we achieve 1.6× smaller energy consumption.
- If the ratio is 10, 2.6× smaller energy consumption is observed.
- For the extreme case of only reads, we can see a 3.82× smaller energy consumption.

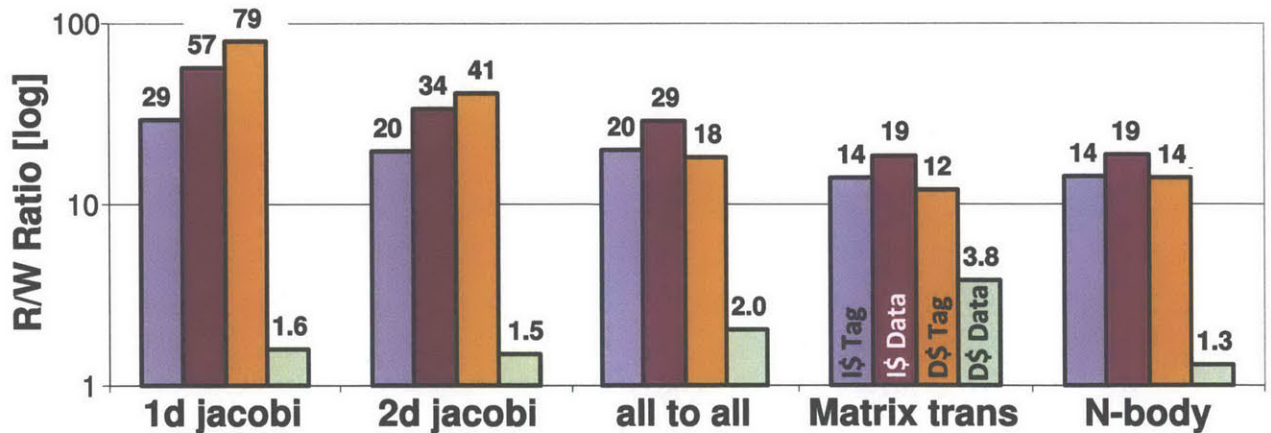


Figure 4-41: Read to write ratios for different benchmarks.

The read to write ratios for different benchmarks are shown in Figure 4-41. This figure is created using the Graphite simulator on a single core processor system with

a 128kB L1-Cache. Five different applications are run in this system and the breakdown of the cache accesses are reported. The figure shows the ratio of the number of total read accesses to write accesses to the L1 cache. Four different columns illustrate the results for: 1- instructor cache tag memory (I\$ Tag), 2-instruction cache data memory (I\$ Data), 3-data cache tag memory (D\$ Tag) and, 4- data cache data memory (D\$ Data). The ratio ranges from 1.3 to 79 between the benchmarks and different memory types. For the instruction cache, the ratio is very high with an average of 32. This is due to the fact that if a particular memory location is referenced, it is likely that the nearby memory location will be referenced in the future. In other words, the instruction caches experience a high spatial locality. According to our test-chip results, this ratio would bring a $3.3\times$ energy savings. On the other hand, the data memory accesses are more random, which decreases the cache hits and requires more lines to be written back from the main memory. The worst case is experienced for the data-cache data bits. For this memory type, an average ratio of 2.0 is observed. Based on the test-chip results, body-biasing would bring 37% energy savings compared to no body-biasing for this ratio.

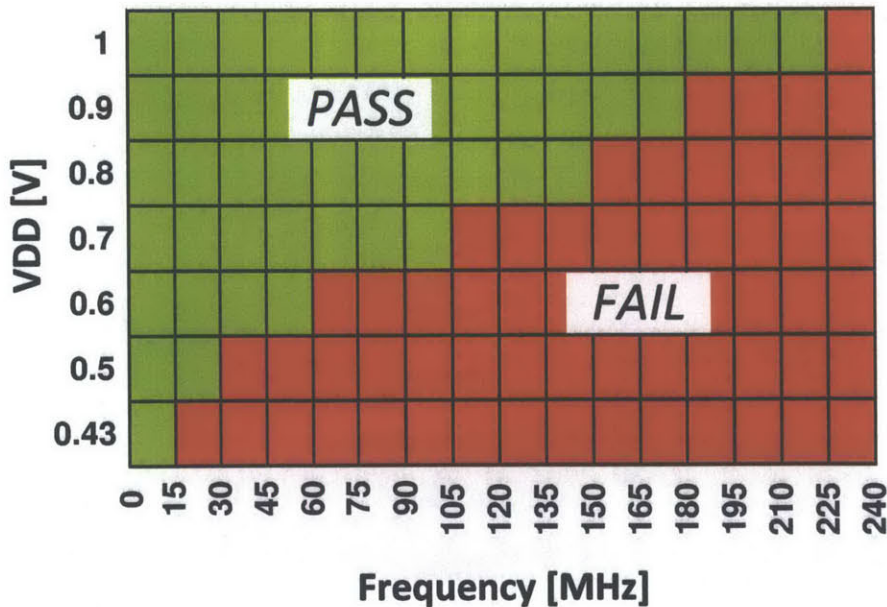


Figure 4-42: Shmoo plot for the 0.5Mb SRAM test-chip designed in 28nm FD-SOI technology.

Figure 4-42 shows the performance vs. V_{DD} shmoo plot of the memory. The

memory achieves its highest performance at 1.0V at 220MHz. The performance scales down to 5MHz at 0.43V.

Table 4.3 shows a comparison of the test-chip with recent SRAM circuits. They are fabricated using scaled technologies of 28nm CMOS and beyond and all use industry sized 6T or 8T bit-cells. Furthermore, all those memories target low-voltage operation and are designed with read and write-assist circuits. As it can be seen from the table, thanks to the memory organization and body-biasing, our memory achieves the lowest minimum voltage operation of 0.43V as opposed to 0.5V [106], 0.6V [56] and 0.7 V [85].

4.3 Summary and Conclusions

SRAM circuits are one of the fundamental building blocks of today's systems and there is a tremendous interest in achieving energy-efficient SRAM design. One of the most effective techniques to achieve energy efficiency is using voltage scaling. However, traditional 6T bit-cell is a ratioed circuit and achieving a reliable operation using the 6T bit-cell at low-voltages is a difficult task. The bit-cell starts experiencing read, write and retention related errors when the voltage is reduced. Furthermore, due to increased variation effects, read-currents are degraded and sensing times are deteriorated.

8T bit-cell is an alternative to the 6T bit-cell since it does not have the read upset problem. Additionally, it increases the options for improving write stability since read and write ports are decoupled from each other. Therefore, in this thesis two SRAM proto-types are designed using 8T bit-cell topology. The first one uses a 65nm LP CMOS technology and it can be voltage scaled from 1.2 V to 0.37 V. The second one is designed using a 0.18 μ m LP CMOS technology and its voltage scalable from 1.8 V to 0.6 V. Both SRAMs use WL boosting to improve write-ability.

The advantages of the 8T bit-cell do not come for free. Firstly, the 8T bit-cell is almost 40% larger compared to its 6T counterpart. Secondly, it does not work with column interleaving, and SAs cannot be shared between multiple columns. This

Table 4.3: Comparison of results with recent publications. *BL-tracked NBL, **suppressed coupling signal for NBL, ***write-recovery-enhancement lower-cell- V_{DD} .

Publication	Technology	Bit-Cell Type & Area	Chip-Area	Capacity	V_{DD} Range	Write-Assist	Read-Assist
	CMOS	μm^2	mm^2	Mb	V		
Sinangil, ISSCC, 2011 [56]	28nm Bulk	6T (0.12)	5.29	0.24	1.0 - 0.6	WL boosting	Short BLs
Kulkarni, ISSCC, 2012 [66]	22nm Bulk	8T (0.238)	NA	1	0.14 improvement	WL boosting	NA
Karl, ISSCC, 2012 [85]	22nm Tri-gate	6T (0.092)	NA	162	1.0 - 0.7	VDD Collapse	WLUD
Chang, ISSCC, 2013 [70]	20nm Tri-gate	6T (0.081)	40.3	112	0.2 improvement	BT-NBL*	PSWL
Chen, ISSCC, 2014 [69]	16nm Tri-gate	6T (0.07)	42.6	128	0.3 improvement	SCS-NBL**, WRE-LCV***	NA
Song, ISSCC, 2014 [106]	14nm Tri-gate	6T (0.064)	75.6	1	1.0 - 0.5	NBL	WLUD
This Work	28nm FD-SOI	6T FW (0.152)	3.24	0.5	1.0 - 0.43	Body-biasing	Short BLs

stresses the area of the SAs significantly for the 8T bit-cell based memories. Smaller sense amplifiers result into degraded SRAM performance at low-voltages since the variation effects are more pronounced. To improve SRAM performance, a SA offset compensation technique is proposed that leverages body-biasing of the input transistors. Although the technique is illustrated for a pseudo-differential sense amplifier, it would also work with a fully differential sense amplifier since it relies on the body voltage. With the requirement of an extra voltage supply for the body, the SA offset is reduced by $2\times$. The calibration process is simple since it requires <5 clock cycles to finish for all the SAs. The area overhead it brings is only 3.5% since the design requires only one latch and a few gates per SA.

Although 8T bit-cell is a viable choice for some designs, it can still be unacceptable due to its larger area and lack of usage with column multiplexing. Therefore, in the second half of this chapter, techniques that make the 6T bit-cell based SRAM design work at lower voltages are investigated.

As transistor scaling reached 45nm and beyond, it has become harder to lower the supply voltage of the SRAMs due of increased transistor variation. Thus, to enable further scaling, new process technologies are introduced. FD-SOI technology is a planar technology and can deliver reduced leakage and variation compared to bulk. Therefore, a test-chip prototype of a 0.5Mb SRAM design using a high-density, $0.152\mu m^2$ 6T bit-cell is designed in a 28nm FD-SOI technology.

An advantage of the FD-SOI technology is the fact that body-biasing is much more effective. Moreover, it enables an extended body-bias range from -3V up to 3V. This provides a new knob in energy efficiency optimization. In the prototype test-chip, extensive body-biasing is used to improve write-ability of the SRAM bit-cell. A flip-well process is used that has NMOS transistors on n-well and PMOS transistors on p-well. This way, the body voltages of the pass-gate transistors are biased to achieve better write-ability. Three techniques are used to reduce the energy consumption due to body-biasing: (1) Stepwise body-biasing, (2) Single-sided body-biasing, and (3) Hierarchical SRAM structure to reduce effective body-capacitor.

To improve read-ability, the memory is organized such that it would have sub-

blocks with 32 bits per column. This improves the read μ/σ from 4.8 to 6.2 at 400mV.

The measurement results of the 28nm FD-SOI SRAM show that this memory is operational down to 430mV. This is the lowest V_{min} reported at a scaled process technology like 28nm FD-SOI using 6T bit-cells. By using a 1.8V body-biasing, a 200 mV reduction in V_{min} is achieved, which, in turn results into energy savings of 37% to $3.8\times$ for a R/W ratio of 2 and higher.

Chapter 5

Conclusions and Future Work

This thesis presents several circuit and design techniques to enable a self-aware processor system. For self-awareness, systems need to be able to observe the important metrics such as energy and performance during runtime. Current systems use performance counters and model-based energy (or power) sensors for observation. These models can come short of representing the dynamically changing operating conditions (e.g. voltage fluctuations, temperature gradients, process corners, aging effects, changing application data) of today's complex computing systems. In this thesis, an energy monitoring technique is developed to measure the absolute energy consumption of a block by introducing minimal area and power overhead. A test-chip prototype is designed to demonstrate the energy monitoring circuit for an SRAM application. Then, the energy monitoring idea is extended to a self-aware processor system. Multiple energy monitors are integrated into the system. Energy measurements can be initiated and results can be read by the software. This way, the software decision engine can optimize the system based on dynamic changes and software targets. This makes it a truly *open* system, where adaptations and measurements are open to software.

Secondly, self-aware systems require reconfigurable circuits or architectures. In this thesis, we propose DVFS and d-cache adaptation for reconfigurability. The d-caches are designed to be reconfigurable to adjust the system for compute- vs. cache-bound applications. To address fast tag invalidation between configuration changes,

custom tag memories are designed to provide single cycle memory reset. To enable efficient voltage conversion, on-chip DC-DC converters are designed; and, energy sensors are embedded into the DC-DC converters. This is the first demonstration of an *open* self-aware processor system that can perform dynamic power/performance optimizations based on actual energy consumption information.

This thesis also focuses on low-voltage SRAM design and performance improvement. 8T bit-cell based memories are demonstrated to perform extensive voltage scaling. However, 8T designs cannot be column-interleaved and their SA area is limited. Therefore, in this thesis a SA offset compensation technique is proposed to achieve input offset reduction. Furthermore, 8T bit-cells bring a significant (40%) area overhead. Therefore, it is important to explore new techniques to achieve low voltage operation with industry sized 6T bit-cells. FD-SOI is a new process, enabling extensive body biasing and decreased transistor variation. Thus, this thesis focuses on low-voltage operation of a 6T bit-cell based SRAM designed in 28nm FD-SOI technology. Assist techniques are proposed to increase cell robustness at low-voltages. To be able to achieve net energy-savings, we propose several methods to decrease energy consumption due to assists (e.g. stepwise and single-sided charging of capacitors, hierarchical structures). A 0.5Mb SRAM using an industry-sized 6T bit-cell is designed in 28nm FD-SOI and is demonstrated to achieve the minimum operating voltage of 0.43V. This is the lowest voltage achieved for a 6T bit-cell design in 28nm technology.

5.1 Summary of Contributions

In this section, we will summarize the key ideas and contributions presented in this thesis.

5.1.1 Energy Monitoring Circuit Demonstration for an SRAM Application

The main contributions are:

- **An energy monitoring circuit that can measure the absolute energy consumption of any block.**
- Implementation of the circuit for an SRAM application in a 65nm LP CMOS test-chip. Measurement results show an accuracy of 10%. The dynamic power overhead is $< 0.1\%$ and the area overhead is 16%.

Some additional important work performed in this section are summarized below:

- Analyzing the effect of energy monitoring circuit challenges such as its effect on SRAM robustness.
- Analyzing the error sources of the energy monitoring operation.
- Investigation of circuits to keep the area and power overhead of the monitoring circuit small.

5.1.2 Self-Aware System That can Reconfigure Itself Based On Actual Energy Measurements

This work is a collaboration of a group of students and the contributions are summarized in Chapter 1.

Self-Aware Processor System

- **Demonstration of a self-aware processor system that is designed to reconfigure itself based on actual energy consumption measurements during run-time. This is an *open* system since the energy numbers and hardware reconfigurability are accessible and controllable by the software decision engine, SEEC.**
- Implementation of a $0.18\mu\text{m}$ CMOS test-chip featuring the circuits and techniques above. The measurement results demonstrate that the energy consumption can be reduced by $5.5\times$ due to voltage scaling and an additional $1.5\times$ due to d-cache reconfigurability.

- Integration of dedicated energy monitoring registers into the core and using pipeline arithmetic blocks for energy calculation.

DC-DC Converters with Embedded Energy Monitors

- **Integration of the energy monitoring circuits into the DC-DC converters.**
- Design of two DC-DC converter circuits to generate output voltages from 1.8V down to 0.6V.
- The filtering capacitor is also used as the storage capacitor of the monitoring circuit. Thus, the energy monitors of this design, do not require extra off-chip components.
- Energy monitors require $<1\%$ area and $<0.1\%$ power overhead.
- An asynchronous boundary to handle communication between different domains.

Reconfigurable and Voltage Scalable Caches

- **Design of a reconfigurable d-cache with set-associativity (1-4 sets) and size (1-4 kB/set) adaptations with a tag invalidation circuit that can perform resetting in a single clock cycle.**

5.1.3 Low-Voltage SRAM Design

8T bit-cell based SRAM

- **A SA offset compensation technique using body-biasing that can achieve $2\times$ smaller offset after calibration. The calibration technique requires <5 clock cycles for calibrating all the SAs.**
- Demonstration of two 8T bit-cell based SRAMs using write-assists: (1) a test-chip prototype in 65nm CMOS with voltage scaling from 1.2 V down to 0.37 V,

and (2) a test-chip prototype in $0.18\mu\text{m}$ CMOS with voltage scaling from 1.8 V down to 0.6 V.

6T bit-cell based SRAM

- **A 0.5Mb SRAM designed in 28nm FD-SOI technology which is measured to operate from 1 V down to 0.43 V. This is the minimum voltage reported in a 6T bit-cell based SRAM in 28nm.**
- Demonstration of cycle-by-cycle, extensive body-biasing (from 0V to 1.8V) to improve write-ability.
- Three techniques to decrease energy consumption due to body-biasing:
 1. Single-sided body-biasing (only from the side where a data zero is written). Compared to body biasing from both sides, we showed that it has almost the same effect on write-ability.
 2. Using step-wise adiabatic charging of the body-capacitor.
 3. Using a smaller body capacitor by partitioning the memory into 32×128 sub-blocks.

5.2 Future Work

This section examines the possibilities of extending the ideas presented in this thesis for future applications.

- *Improving the energy sensing circuit.*

This thesis demonstrated an energy sensing circuit for monitoring the energy consumption of (1) an SRAM application and, (2) multiple blocks within a processor system. However, the energy sensing circuit can still be improved. With the current design, the accuracy of the sensor depends on the comparator offset significantly. To improve the accuracy, even more sophisticated offset

compensation techniques can be used for the comparator. Secondly, for absolute energy calculation, the value of the capacitor is required. In this design, we assumed that the capacitor value is known. As mentioned, measuring the actual capacitance can be necessary to improve the accuracy of the actual energy consumption. This can be performed during a calibration phase prior to operation.

- *Self-Awareness for Multi-Core Processors.*

In recent processor systems, power control units (PCU) that estimate power consumption based on energy models are used for self-adaptation. Monitoring the actual energy consumption of different blocks within a processor can capture the dynamic effects of the system better. Therefore, the usage of energy sensing circuit can be extended to multi-core processor designs. Similarly, for multi-core systems, another knob for reconfigurability can be turning some of the cores ON and OFF depending on the system condition. Alternatively, specialized, low-power cores can be used for manipulating the states. For the Angstrom Processor, *Partner Cores* were conceptually proposed for this reason [107]. The next step for the self-aware microprocessor design is to extending it to multi-core systems and analyzing its trade-offs in an actual hardware system.

- *Sense Amplifier Offset Compensation Using Multiple Supply Voltages of Body Biasing.*

The offset compensation technique proposed in this thesis uses one supply voltage for the body voltage and improves the input offset compensation by $2\times$. However, theoretically, with the addition of every extra voltage supply, the offset can be reduced by an additional $2\times$. Therefore, another technique that can be investigated is using multiple sources for reducing the sense amplifier offset. Secondly, since our design uses the body voltage for offset compensation, it can also be used in differential sense amplifiers. Therefore, the concept can be investigated in a differential sense amplifier. Thirdly, the current design assumes that the reference voltage and the body voltage are generated outside.

The next step can be integrating the voltage converters for that on-chip. Furthermore, a control loop can be used to determine to pick the reference voltage automatically.

- *New 6T bit-cell Layout with Horizontal n-well stripes to Achieve a Lower Body Capacitance.*

In this thesis, we used extensive body biasing capability of FD-SOI to improve write-ability. We showed our cycle-by-cycle body-biasing idea for an SRAM design that uses industry sized 6T bit-cell. We introduced minimal changes to the SRAM array structure. In other words, we demonstrated the idea for minimum changes to the existing SRAM design so that it can easily be adopted. However, in the conventional bit-cell structure, n-wells are routed column-wise. Thus, while we bias the body capacitor of the accessed cell during a write operation, the body capacitor of the non-accessed rows within the same sub-block are also biased. This is not ideal since it increases the energy consumption due to body biasing significantly. Alternatively, a new layout of the 6T bit-cell can be designed with row-wise n-wells. This way, the body capacitor can be made extensively smaller. For the new bit-cell to be able to be an alternative to the conventional 6T bit-cell, it should not introduce a large area overhead and still needs to preserve the high-density structure of the 6T bit-cell.

Appendix A

Yield Calculation of SRAMs

SRAMs account for a significant portion of the total transistor count in many systems. In a typical processor system, this number is around 80% [108]. Thus, yield management of these SRAMs play a crucial role in ensuring design success.

This appendix demonstrates a simple analysis to relate SRAM bit-cell failure rate to its yield requirement. The limits that are provided here can be relaxed by inserting redundant rows and columns, Error Correction Coding (ECC) or increasing the size of the cells [109, 62].

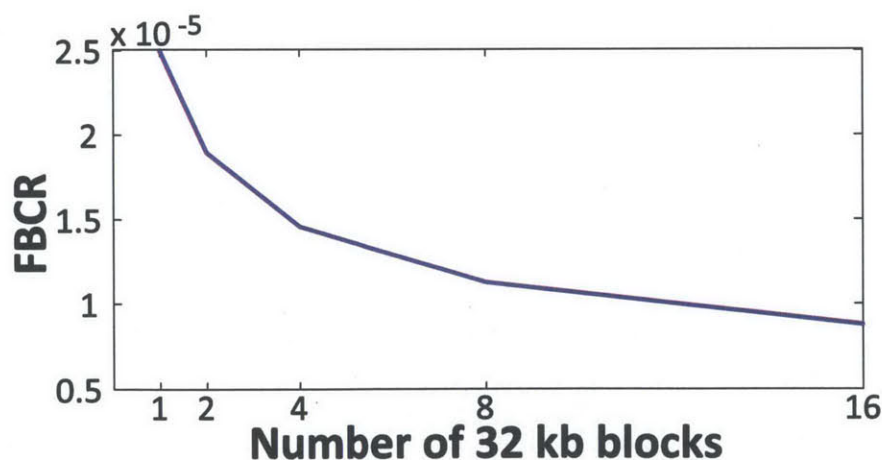


Figure A-1: FBR vs. the memory size for 95% yield

Fail-bit ratio (FBR) is the ratio of failing bit-cells to the total number of bit-cells in the SRAM array. This ratio can be used to determine SRAM yield for a fixed memory size [88]. Figure A-1 shows FBR that can be tolerated for different memory

sizes. This analysis assumes that a memory block of 32 kb is replicated to generate a larger memory sizes and every 32 kb memory block uses 2 redundant rows. As shown in the figure, to guarantee same yield, a smaller FBR constraint has to be selected as the memory size increases.

Appendix B

SRAM Assist Techniques

This section gives a summary of various bit-cell assist techniques used in recent SRAM designs. Table B.1 gives examples for write-ability improvement; whereas, Table B.2 gives examples for read-ability improvement.

Table B.1: Various Bit-cell Assist Techniques for Enhancing **Write-ability**.

Work	Technique	Description
Karl, ISSCC, 12, [85]	Dynamic modulation of VDD	Actively drive selected vertical VDD to VSS
Kulkarni, ISSCC, 12, [66]	Dynamic modulation of VDD	Self-induced VDD collapse (SIC)
Chen, ISSCC, 14, [69]	Dynamic modulation of BL	Suppressed-coupling-signal NBL
Chang, ISSCC, 14, [70]	Dynamic modulation of BL	Bit-line length tracked NBL
Sinangil, VLSI, 14, [73]; Ickes, ESSCIRC, 11, [37]	Static WWL boosting	Use VWL > VDD for 8T
Yamaoka, ESSCIRC, 07, [59]	Body-biasing	Recenter the balance between write and read with variability

Table B.2: Various Bit-cell Assist Techniques for Enhancing **Read-ability**.

Work	Technique	Description
Nho, ISSCC, 10, [52]	Dynamic modulation of WL	Adaptive WLUD by tracking PVT
Khellah, VLSI, 06, [54]	Dynamic modulation of WL	During read, cut-off WL pulse before unstable bit-flip completes
Pilo, JSSC, 12, [55]	Dynamic modulation of BL	Precharge bitlines to a lower voltage
Sinangil, ISSCC, 11, [56]	Memory organization	Use a small BL capacitance
Verma, ISSCC, 07 [62]	Dynamic modulation of horizontally routed VSS	For 8T, raise footer of the read-buffer for the unselected rows to eliminate bitline off current
Yabuuchi, VLSI, 09 [110]	Dynamic modulation of vertically routed VSS	Negative VSS using a charge pump

Appendix C

LEON3 Processor

This appendix gives a brief description of the original LEON3 (vanilla LEON3) architecture. Some of the important specifications of the vanilla LEON3 are given in Table C.1.

Table C.1: Vanilla LEON3 processor core features.

Feature	LEON3
Architecture	SPARC V8 instruction set
Bit Width	32
Pipeline Stages	7
FPU	Optional
Branch Prediction	Static
Cache Configuration	1-4 ways, 1-256 kbytes/way
Bus Interface	AMBA-2.0 AHB bus
DRAM Controller	16/32/64-bit DDR/DDR2 controllers
MMU Page Size	Fixed
Number of Gates	30K
Debug Support	GRMON

The LEON3 is a synthesizable VHDL model of a 32-bit processor which is compliant with the Scalable Processor ARChitecture (SPARC) V8 RISC [111]. It supports the IP core of Gaisler Research IP library (GRLIB), which is an integrated set of reusable IP cores. LEON3 is particularly suitable for SoC designs like the self-aware processor since it is highly configurable, and relatively simple. It is a simple design since it performs no out-of-order execution (OOE), has a static branch prediction,

and can be implemented using approximately 30K gates.

This processor and its derivatives have been used both in professional and in academic research applications. One example is the Berkeley Emulation Engine 2 (BEE2) [112] which is a multi-chip FPGA board. Another example is the processor in [113] which is designed using an enhanced LEON3 core with a low-power management unit. The self-aware processor system that is presented in this thesis [73] is also designed based on the LEON3 processor.

The LEON3 pipeline includes functionality to allow non-intrusive debugging on the target hardware. To aid software debugging, watchpoints or breakpoints can be enabled. Through the debug support interface, full access to all processor registers and caches is provided. It also enables single stepping and hardware breakpoint/watchpoint control. Additionally, an internal trace buffer monitors and stores the executed instructions, which can be later read out over the debug interface. GRMON is a debug monitor for the vanilla LEON3 debug support unit (DSU), providing a non-intrusive debug environment on the real target hardware. The DSU can be controlled through the AHB. These features are very desirable to debug a complicated system design like the self-aware processor and adopted by it.

LEON3 comes with a highly configurable cache system, consisting of a separate instruction-cache (*i-cache*) and data-cache (*d-cache*). Both caches can be configured with 1-4 sets, 1kB to 256kB per set, and 16B or 32B per line. This feature is very important for the self-aware processor since in this system, d-caches are designed to be reconfigurable.

The LEON3 core is centered around the AMBA Advanced High-speed Bus (AHB), to which the core and other blocks are connected. It implements a 7-stage integer pipeline with optional multiply and divide units. Similarly, vanilla LEON3 support the MAC instructions of SPARC V8 instruction set. Those are specifically useful for DSP algorithms.

The vanilla LEON3 processor core has an optional power-down mode, which halts the pipeline and caches. This is particularly useful for systems with multi-core designs. The vanilla LEON3 processor also provides interfaces for a co-processor and it has

the ability to work with an floating-point unit (FPU) block.

Appendix D

Switching Regulators

There are several power supply voltage regulator topologies. The most popular three are the linear regulators, switched capacitor converters and switching regulators. In this appendix, we will summarize the switched capacitor circuits.

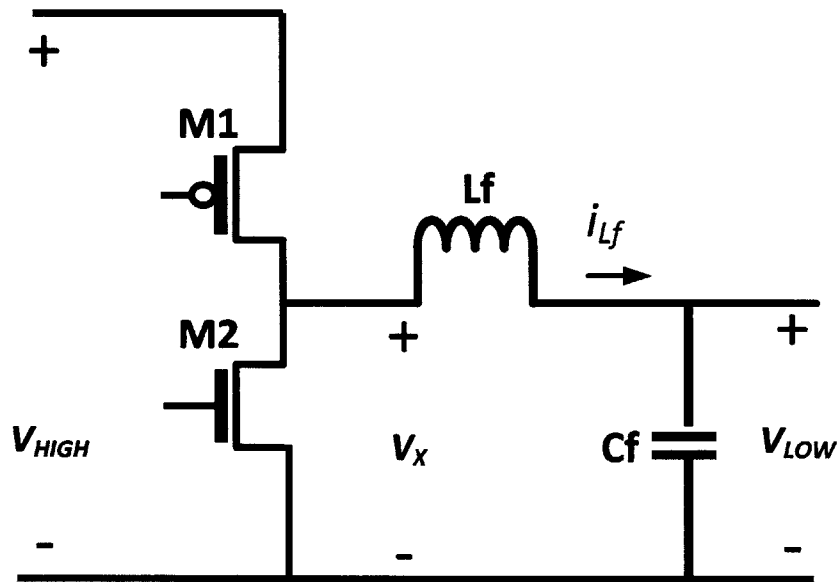


Figure D-1: A typical buck converter.

The unique advantage of the switching regulators lies in their ability to convert a given supply voltage with a known voltage range to virtually any given desired output voltage, with no first order limitations on the efficiency. This is true regardless of whether the supply voltage is higher or lower or with different polarity than the input voltage. The basic components of a switching regulator are inductor and capacitor,

which are, ideally, reactive elements that dissipate no power. Similarly, the switches are ideally either ON or OFF which means that either the current or the voltage across them is equal to zero. Thus, the power consumption across them is ideally equal to zero. Under the ideal components assumption, the switching regulators are 100% efficient.

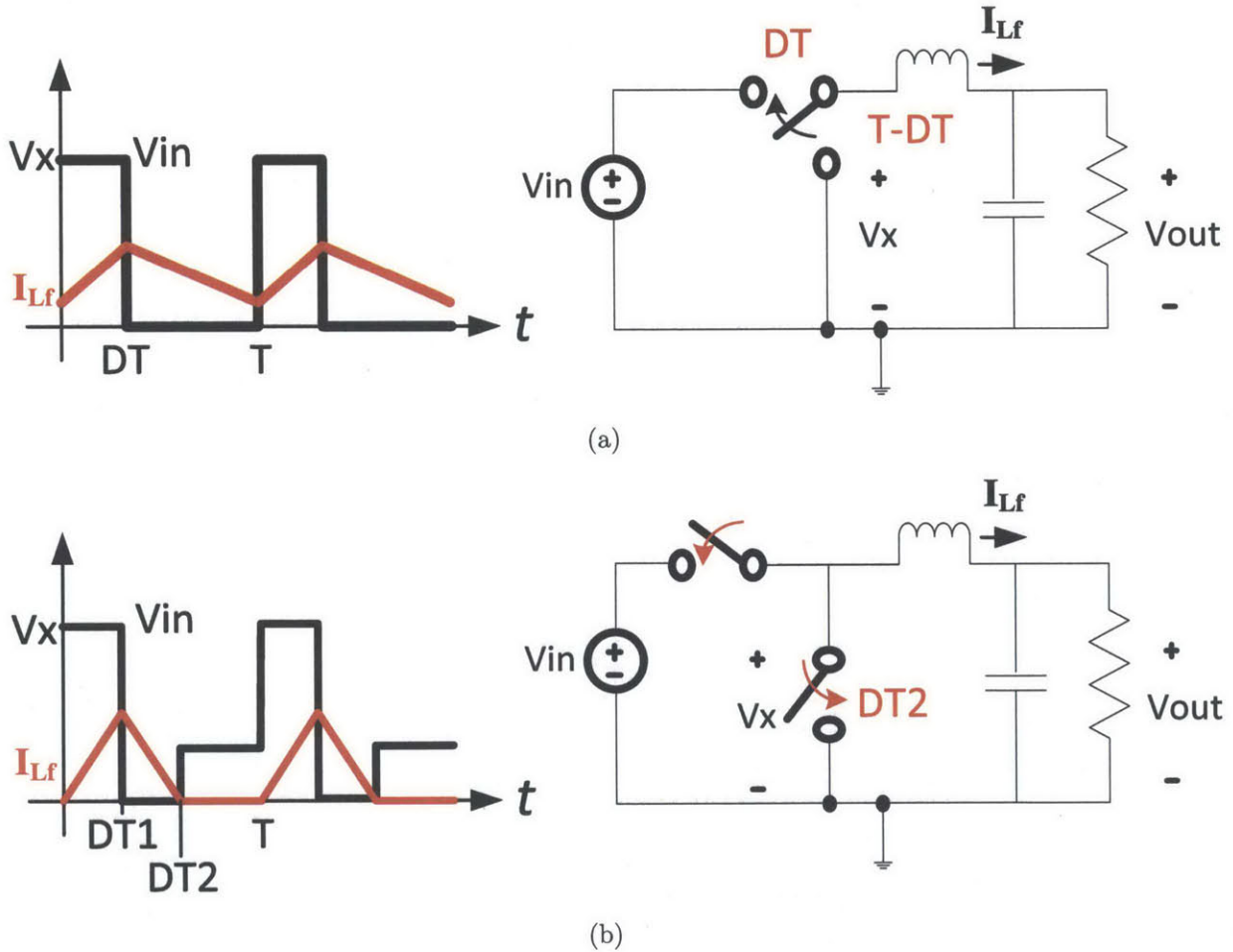


Figure D-2: Operation of a buck converter in (a) PWM mode control, and (b) in PFM mode control.

Depending on the topology of the switching regulator, it can be used either to step-down or up a voltage. The popular topology to step-down a voltage is the buck-converter. Similarly, the popular topology to step-up a voltage is the boost-converter. The typical implementation of a buck converter is shown in Figure D-1 which consists of the power transistors M1 and M2 and the filter elements Lf and Cf. It is used to

step-down the voltage from V_{HIGH} to V_{LOW} .

Switching of M1 and M2 power transistors can be controlled in two different schemes in a buck converter. The first one is the Pulse Width Modulation (PWM) mode. In this method, the power transistors generate a chopped version of V_{IN} at V_X node as shown in Figure D-2 (a). The average of V_X is always lower than the average value of V_{IN} . Then, this chopped voltage is low-pass-filtered to generate a low-ripple output voltage of $V_{OUT} = V_{IN}/D$. In this controlling scheme, the transistors are switched at a constant frequency which is equal to $1/T$.

A different mode of operation is called the Pulse Frequency Modulation (PFM) mode where the transistors are switched ON only when the output voltage falls below the desired reference voltage V_{REF} . Operation of this mode can be shown in Figure D-2 (b). Every switching cycle, PMOS is turned on first ramping up the inductor current. Then it is turned OFF after a specific time of T_{PMOS} and NMOS is turned ON and remains ON ramping down the inductor current till it hits zero. The time for which the NMOS has to be kept on is determined by V_{LOW} , V_{HIGH} and T_{PMOS} . In this scheme, the frequency of operation depends on T_{PMOS} , T_{NMOS} , the leakage current of the loading circuit and the ripple voltage at the output.

In the design a PFM mode buck converter is used for voltage conversion due to two reasons. Firstly, in the PFM mode, the current through the inductor is not continuous and the transistors are OFF for a major part of the clock cycle. This is more convenient for the energy monitoring operation since energy monitoring requires the transistors to be shut off. Otherwise, if a PWM mode converter was used with the energy sensors, the stored charge on the inductor would need to be wasted every energy sensing cycle. Secondly, the PFM mode is more suitable for low-voltage operation and it provides higher efficiencies.

There are different sources of power loss in the buck converters. In this section, we will explain the main sources of inefficiency for buck converters with PFM mode modulation. The main sources can be classified as conduction losses, switching losses, and timing losses.

D.0.1 Conduction Loss

The conduction loss is due to the current flowing through the components such as power transistors, filtering capacitor, inductors, and wires.

Since in the PFM mode the current through the inductor is not continuous, the transistors are OFF for the majority of the time, it is more convenient to talk about the efficiency by energy loss. The energy dissipated in one charge cycle due to conduction is given by:

$$E_{cond} = \int_0^T i(t)^2 R dt. \quad (D.1)$$

This equation has two parts. During the first part, the PMOS transistor (M1 in Figure D-1) is ON and the resistor is the ON resistance of the PMOS transistor. During this stage, the transistor is operating in its triode region and its current is proportional to the inverse of its width (W_{PMOS}). The time for the integral for this period is the ON time of the PMOS. This is shown as DT1 in Figure D-2 (b). For the second part of the equation, the NMOS transistor (M2) is ON and its resistance is related to $1/W_{NMOS}$. Similarly, the time period for this time is DT2.

D.0.2 Switching Loss

Since the filtering transistors of the buck converter need to deliver the load current and since they have conduction loss associated with them, they need to be designed considerably large. However, due to that reason, their gate capacitances are tend to large and charging up the gate capacitances of those transistors requires a considerable energy consumption. The energy consumed for switching in one cycle of operation of the PFM mode converter is given by E_{sw} . This energy is proportional to the width of the switching transistors and hence can be written as:

$$E_{sw} = E_{gate} \times W \quad (D.2)$$

Here, E_{gate} is a constant and need to be determined per transistor. Similar to the conduction loss, both NMOS and PMOS transistors need to be considered for this equation.

D.0.3 Timing Losses

The timing losses happen due to non-idealities of switching times for the filtering transistors. There are different mechanisms to the timing losses. For instance, if there is overlap between period DT1 (when PMOS is ON) and DT2 (when NMOS is ON); both transistors can be on at the same time. A large short circuit current can flow. This is called the *short circuit loss*. On the other extreme, if there is a very long dead time between DT1 and DT2, NMOS body-diode can start to conduct.

D.0.4 Leakage Loss

Although not extremely important for technologies such as $0.18\mu\text{m}$, leakage current of the filtering transistors might be significant for scaled technologies. For the PFM mode DC-DC converter, the energy lost due to leakage of power transistors in one switching cycle is given by $E_{leak} = Vin \times T$.

There are also losses associated with the control circuitry and pulse generation circuitry that might be significant based on the design.

D.0.5 Efficiency of the PFM Mode DC-DC Converter

Considering all the losses, the efficiency of the PFM mode converter can be summarized as:

$$\nu = \frac{E_{load}}{E_{load} + E_{losses}} \quad (\text{D.3})$$

where $E_{losses} = E_{cond} + E_{sw} + E_{timing} + E_{leak} + E_{other}$.

Bibliography

- [1] R. Riedlinger, R. Bhatia, L. Biro, B. Bowhill, E. Fetzer, P. Gronowski, and T. Grutkowski, "A 32nm 3.1 billion transistor 12-wide-issue itanium processor for mission-critical servers," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2011, pp. 84–86.
- [2] E. Fluhr, J. Friedrich, D. Dreps, V. Zyuban, G. Still, C. Gonzalez, A. Hall, D. Hogenmiller, F. Malgioglio, R. Nett, J. Paredes, J. Pille, D. Plass, R. Puri, P. Restle, D. Shan, K. Stawiasz, Z. Deniz, D. Wendel, and M. Ziegler, "Power8™: A 12-core server-class processor in 22nm soi with 7.6tb/s off-chip bandwidth," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2014, pp. 96–97.
- [3] "ISSCC Technology Trends 2013." [Online]. Available: <http://isscc.org/doc/2013/2013.trends.pdf>.
- [4] M. H. Abu-Rahma and M. Anis, *Nanometer Variation-Tolerant SRAM*. Springer Science and Business Media, 2013.
- [5] M. Qazi, M. Sinangil, and A. Chandrakasan, "Challenges and directions for low-voltage sram," *Design Test of Computers, IEEE*, vol. 28, no. 1, pp. 32–43, Jan 2011.
- [6] E. Rotem, A. Naveh, D. Rajwan, A. Ananthakrishnan, and E. Weissmann, "Power-management architecture of the intel microarchitecture code-named sandy bridge," *IEEE Micro*, vol. 32, no. 2, pp. 20–27, March 2012.
- [7] R. Ranica, N. Planes, O. Weber, O. Thomas, S. Haendler, D. Noblet, D. Croain, C. Gardin, and F. Amaund, "Fdsoi processdesign full solutions for ultra low leakage, high speed and low voltage srams," in *Symp. on VLSI Technology (VLSI) Dig. Tech. Papers*, 2013.
- [8] "Innovation and Technology: FD-SOI." [Online]. Available: http://www.st.com/web/en/about_st/learn_fdsoi.html.
- [9] G. E. Moore, "Cramming More Components onto Integrated Circuits," *Electronics*, vol. 38, no. 8, pp. 114–117, Apr. 1965.

- [10] Y. K. Ramadass, “Energy Processing Circuits for Low-Power Applications,” Ph.D. dissertation, Massachusetts Institute of Technology, Electrical Engineering and Computer Science, 2009.
- [11] “Panasonic Develops New Higher-Capacity 18650 Li-Ion Cells; Application of Silicon-based Alloy in Anode.” [Online]. Available: greencarcongress.com.
- [12] “Battery Statistics.” [Online]. Available: batteryuniversity.com.
- [13] A. Chandrakasan, D. Daly, J. Kwong, and Y. Ramadass, “Next generation micro-power systems,” in *Symp. on VLSI Circuits (VLSI) Dig. Tech. Papers*, June 2008, pp. 2–5.
- [14] J. Warnock, Y. H. Chan, H. Harrer, D. Rude, R. Puri, S. Carey, G. Salem, G. Mayer, Y. Chan, M. Mayo, A. Jatkowski, G. Strevig, L. Sigal, A. Datta, A. Gattiker, A. Bansal, D. Malone, T. Strach, W. Huajun, P. Mak, C. Shum, D. Plass, and C. Webb, “5.5GHz system z microprocessor and multi-chip module,” *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, pp. 46–47, Feb. 2013.
- [15] K. Gillespie, H. R. Fair, C. Henrion, R. Jotwani, S. Kosonocky, R. S. Orefice, D. A. Priore, J. White, and K. Wilcox, “Streamroller: An x86-64 Core Implemented in 28nm Bulk CMOS,” *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, pp. 104–105, Feb. 2014.
- [16] W. Hu, Y. Zhang, L. Yang, B. Fan, Y. Chen, S. Zhong, H. Wang, Z. Qi, P. Wang, X. Gao, X. Yang, B. Xiao, H. Wang, Z. Yang, L. Yang, and S. Chen, “Godson-3b1500: A 32nm 1.35ghz 40w 172.8gflops 8-core processor,” in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2013, pp. 54–55.
- [17] S. Choi and D. Yeung, “Learning-based smt processor resource distribution via hill-climbing,” in *Computer Architecture, 2006. ISCA '06. 33rd International Symposium on*, 2006, pp. 239–251.
- [18] D. Albonesi, R. Balasubramonian, S. Ddropsbo, S. Dwarkadas, E. Friedman, M. Huang, V. Kursun, G. Magklis, M. Scott, G. Semeraro, P. Bose, A. Buyuktosunoglu, P. Cook, and S. Schuster, “Dynamically tuning processor resources with adaptive processing,” *IEEE Computer*, vol. 36, no. 12, pp. 49–58, Dec 2003.
- [19] C. Dubach, T. Jones, E. Bonilla, and M. O’Boyle, “A predictive model for dynamic microarchitectural adaptivity control,” in *Microarchitecture (MICRO), 2010 43rd Annual IEEE/ACM International Symposium on*, Dec 2010, pp. 485–496.
- [20] H. Hoffmann, J. Holt, G. Kurian, E. Lau, M. Maggio, J. E. Miller, S. M. Neuman, M. Sinangil, Y. Sinangil, A. Agarwal, A. P. Chandrakasan, and S. Devadas, “Self-aware Computing in the Angstrom Processor,” *ACM/IEEE Design Automation Conf. (DAC) Dig. Tech. Papers*, pp. 259–264, 2012.

- [21] Y. Ramadass and A. Chandrakasan, "Minimum energy tracking loop with embedded dc-dc converter delivering voltages down to 250mv in 65nm cmos," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2007, pp. 64–587.
- [22] M. Yuffe, E. Knoll, M. Mehalel, J. Shor, and T. Kurts, "A Fully Integrated Multi-CPU, GPU and Memory Controller 32nm Processor," *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, pp. 264–266, Feb. 2011.
- [23] S. Damaraju, V. George, S. Jahagirdar, T. Khondker, R. Milstrey, S. Sarkar, S. Siers, I. Stoloro, and A. Subbiah, "A 22nm IA multi-CPU and GPU System-on-Chip," *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, pp. 56–57, Feb. 2012.
- [24] S. Rusu, H. Muljono, D. Ayers, S. Tam, W. Chen, A. Martin, S. Li, S. Vora, R. Varada, and E. Wang, "Ivytown: A 22nm 15-core enterprise xeon processor family," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2014, pp. 102–103.
- [25] R. Kan, T. Tanaka, G. Sugizaki, R. Nishiyama, S. Sakabayashi, Y. Koyanagi, R. Iwatsuki, K. Hayasaka, T. Uemura, G. Ito, Y. Ozeki, H. Adachi, K. Furuya, and T. Motokurumada, "A 10th generation 16-core sparc64 processor for mission-critical unix server," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2013, pp. 60–61.
- [26] A. Chandrakasan, S. Sheng, and R. Brodersen, "Low-Power CMOS Digital Design," *IEEE J. Solid-State Circuits*, vol. 27, no. 4, pp. 473–484, Apr. 1992.
- [27] D. F. Finchelstein, V. Sze, , M. E. Sinangil, Y. Koken, and A. P. Chandrakasan, "A low-power 0.7V H.264 720 video decoder," in *IEEE Asian Solid State Circuits Conference*, Nov. 2008, pp. 173–176.
- [28] J. Kephart and D. Chess, "The vision of autonomic computing," *IEEE Computer*, vol. 36, no. 1, pp. 41–50, Jan 2003.
- [29] R. Laddaga, "Creating robust software through self-adaptation," *Intelligent Systems and their Applications, IEEE*, vol. 14, no. 3, pp. 26–29, May 1999.
- [30] M. Salehie and L. Tahvildari, "Self-adaptive software: Landscape and research challenges," *ACM Trans. Auton. Adapt. Syst.*, vol. 4, no. 2, May 2009.
- [31] R. Bitirgen, E. Ipek, and J. Martinez, "Coordinated management of multiple interacting resources in chip multiprocessors: A machine learning approach," in *IEEE Micro*, Nov 2008, pp. 318–329.
- [32] H. Hoffmann, M. Maggio, M. D, A. Leva, A. Agarwal, H. Hoffmann, M. Maggio, M. D. Santambrogio, A. Leva, and A. Agarwal, "Seec: A framework for self-aware computing," Tech. Rep., 2010.

- [33] P. Dutta, M. Feldmeier, J. Paradiso, and D. Culler, "Energy Metering for Free: Augmenting Switching Regulators for Real-Time Monitoring," in *Information Processing in Sensor Networks (IPSN)*, 2008, pp. 283–294.
- [34] B. H. Calhoun and A. Chandrakasan, "Characterizing and Modeling Minimum Energy Operation for Subthreshold Circuits," in *Int. Symp. on Low-Power Elec. and Design (ISLPED) Dig. Tech. Papers*, 2004, pp. 90–95.
- [35] A. Wang, A. Chandrakasan, and S. Kosonocky, "Optimal Supply and Threshold Scaling for Sub-threshold CMOS Circuits," in *IEEE Computer Society Annual Symposium on VLSI*, Apr. 2002, pp. 7–11.
- [36] B. Calhoun, A. Wang, and A. Chandrakasan, "Modeling and Sizing for Minimum Energy Operation in Subthreshold Circuits," *IEEE J. Solid-State Circuits*, vol. 40, no. 9, pp. 1778–1786, Sep. 2005.
- [37] N. Ickes, Y. Sinangil, F. Pappalardo, E. Guidetti, and A. P. Chandrakasan, "A 10pJ/cycle ultra-low-voltage 32-bit microprocessor system-on-chip," *IEEE European Solid-State Circuits Conf. (ESSCIRC) Dig. Tech. Papers*, pp. 159–162, 2011.
- [38] J. M. Rabaey, A. Chandrakasan, and B. Nikolic, *Digital Integrated Circuits: A Design Perspective*, 2nd ed. Prentice Hall, 2003.
- [39] M. Sinangil, N. Verma, and A. Chandrakasan, "A Reconfigurable 8T Ultra-Dynamic Voltage Scalable U-DVS SRAM in 65 nm CMOS," *Solid-State Circuits, IEEE Journal of*, vol. 44, no. 11, pp. 3163–3173, Nov. 2009.
- [40] H. Yamauchi, "Embedded sram circuit design technologies for a 45nm and beyond," in *ASIC, 2007. ASICON '07. 7th International Conference on*, Oct 2007, pp. 1028–1033.
- [41] A. Chandrakasan, W. Bowhill, and F. Fox, *High-Performance Microprocessor Circuits*. Wiley-IEEE Press, Piscataway, 2000, 2000.
- [42] B. Wong, A. Mittal, Y. Cao, and G. Starr, *Nano-CMOS Circuit and Physical Design*. Wiley-Interscience, New York, 2004, 2004.
- [43] B. Kiyoo Itoh, M. Horiguchi, and M. Yamaoka, "Low-voltage limitations of memory-rich nano-scale cmos lsis," in *Solid State Device Research Conference, 2007. ESSDERC 2007. 37th European*, Sept 2007, pp. 68–75.
- [44] R. Baumann, "The impact of technology scaling on soft error rate performance and limits to the efficacy of error correction," in *Electron Devices Meeting, 2002. IEDM '02. International*, Dec 2002, pp. 329–332.
- [45] M. Yamaoka, N. Maeda, Y. Shinozaki, Y. Shimazaki, K. Nii, S. Shimada, K. Yanagisawa, and T. Kawahara, "Low-Power Embedded SRAM Modules with Expanded Margins for Writing," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2005, pp. 480–481.

- [46] M. Yabuuchi, K. Nii, Y. Tsukamoto, S. Ohbayashi, S. Imaoka, H. Makino, Y. Yamagami, S. Ishikura, T. Terano, T. Oashi, K. Hashimoto, A. Sebe, S. Okazaki, K. Satomi, H. Akamatsu, and H. Shinohara, "A 45nm low-standby-power embedded sram with improved immunity against process and temperature variations," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2007, pp. 326–606.
- [47] R. Heald and P. Wang, "Variability in sub-100nm sram designs," in *IEEE Int. Conf. on CAD (ICCAD) Dig. Tech. Papers*, Nov 2004, pp. 347–352.
- [48] E. Seevinck, F. List, and J. Lohstroh, "Static Noise Margin Analysis of MOS SRAM Cells," *IEEE J. Solid-State Circuits*, vol. SC-22, no. 5, pp. 748–754, Oct. 1987.
- [49] K. Nii, M. Yabuuchi, Y. Tsukamoto, S. Ohbayashi, S. Imaoka, H. Makino, Y. Yamagami, S. Ishikura, T. Terano, T. Oashi, K. Hashimoto, A. Sebe, S. Okazaki, K. Satomi, H. Akamatsu, and H. Shinohara, "A 45-nm bulk cmos embedded sram with improved immunity against process and temperature variations," *IEEE J. Solid-State Circuits*, vol. 43, no. 1, pp. 180–191, Jan 2008.
- [50] S. Ohbayashi, M. Yabuuchi, K. Nii, Y. Tsukamoto, S. Imaoka, Y. Oda, T. Yoshihara, M. Igarashi, M. Takeuchi, H. Kawashima, Y. Yamaguchi, K. Tsukamoto, M. Inuishi, H. Makino, K. Ishibashi, and H. Shinohara, "A 65-nm soc embedded 6t-sram designed for manufacturability with read and write operation stabilizing circuits," *Solid-State Circuits, IEEE Journal of*, vol. 42, no. 4, pp. 820–829, April 2007.
- [51] O. Hirabayashi, A. Kawasumi, A. Suzuki, Y. Takeyama, K. Kushida, T. Sasaki, A. Katayama, G. Fukano, Y. Fujimura, T. Nakazato, Y. Shizuki, N. Kushiyama, and T. Yabe, "A process-variation-tolerant dual-power-supply sram with 0.179m² cell in 40nm cmos using level-programmable wordline driver," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2009, pp. 458–459,459a.
- [52] H. Nho, P. Kolar, F. Hamzaoglu, Y. Wang, E. Karl, Y. Ng, U. Bhattacharya, and K. Zhang, "A 32nm high-k metal gate sram with adaptive dynamic stability enhancement for low-voltage operation," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2010, pp. 346–347.
- [53] A. Bhavnagarwala, S. Kosonocky, C. Radens, Y. Chan, K. Stawiasz, U. Srinivasan, S. P. Kowalczyk, and M. Ziegler, "A sub-600-mv, fluctuation tolerant 65-nm cmos sram array with dynamic cell biasing," *IEEE J. Solid-State Circuits*, vol. 43, no. 4, pp. 946–955, April 2008.
- [54] M. Khellah, Y. Ye, N. S. Kim, D. Somasekhar, G. Pandya, A. Farhang, K. Zhang, C. Webb, and V. De, "Wordline and bitline pulsing schemes for improving sram cell stability in low-vcc 65nm cmos designs," in *Symp. on VLSI Circuits (VLSI) Dig. Tech. Papers*, 2006, pp. 9–10.

- [55] H. Pilo, I. Arsovski, K. Batson, G. Braceras, J. Gabric, R. Houle, S. Lamphier, C. Radens, and A. Seferagic, "A 64 mb sram in 32 nm high-k metal-gate soi technology with 0.7 v operation enabled by stability, write-ability and readability enhancements," *IEEE J. Solid-State Circuits*, vol. 47, no. 1, pp. 97–106, Jan 2012.
- [56] M. Sinangil, H. Mair, and A. Chandrakasan, "A 28nm high-density 6T SRAM with optimized peripheral-assist circuits for operation down to 0.6V," in *Solid-State Circuits Conference Digest of Technical Papers (ISSCC), 2011 IEEE International*, Feb. 2011, pp. 260–262.
- [57] M. Yamaoka, K. Osada, and T. Kawahara, "A cell-activation-time controlled sram for low-voltage operation in dvfs socs using dynamic stability analysis," in *IEEE European Solid-State Circuits Conf. (ESSCIRC) Dig. Tech. Papers*, Sept 2008, pp. 286–289.
- [58] B. Amrutur and M. Horowitz, "Speed and power scaling of sram's," *IEEE J. Solid-State Circuits*, vol. 35, no. 2, pp. 175–185, Feb 2000.
- [59] M. Yamaoka and T. Kawahara, "Operating-margin-improved sram with column-at-a-time body-bias control technique," in *IEEE European Solid-State Circuits Conf. (ESSCIRC) Dig. Tech. Papers*, Sept 2007, pp. 396–399.
- [60] S. Cosemans, W. Dehaene, and F. Catthoor, "A 3.6pJ/access 480MHz, 128Kbit on-chip SRAM with 850MHz boost mode in 90nm CMOS with tunable sense amplifiers to cope with variability," in *IEEE European Solid-State Circuits Conf. (ESSCIRC) Dig. Tech. Papers*, Sep. 2008, pp. 278–281.
- [61] M. Sinangil, N. Verma, and A. Chandrakasan, "A 45nm 0.5V 8T column-interleaved SRAM with on-chip reference selection loop for sense-amplifier," in *Solid-State Circuits Conference, 2009. A-SSCC 2009. IEEE Asian*, Nov. 2009, pp. 225–228.
- [62] N. Verma and A. Chandrakasan, "A 65nm 8T Sub-Vt SRAM Employing Sense-Amplifier Redundancy," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2006, pp. 328–329.
- [63] S. O. Toh, Z. Guo, T.-J. Liu, and B. Nikolic, "Characterization of Dynamic SRAM Stability in 45 nm CMOS," *Solid-State Circuits, IEEE Journal of*, vol. 46, no. 11, pp. 2702–2712, Nov. 2011.
- [64] K. Kushida and et al., "A 0.7V single-supply SRAM with $0.495\mu\text{m}^2$ cell in 65nm technology utilizing self-write-back sense amplifier and cascaded bit line scheme," in *Symp. on VLSI Circuits (VLSI) Dig. Tech. Papers*, Jun. 2008, pp. 46–47.
- [65] M. Yamaoka, N. Maeda, Y. Shinozaki, Y. Shimazaki, K. Nii, S. Shimada, K. Yanagisawa, and T. Kawahara, "90-nm process-variation adaptive embedded

- sram modules with power-line-floating write technique,” *IEEE J. Solid-State Circuits*, vol. 41, no. 3, pp. 705–711, March 2006.
- [66] J. Kulkarni, B. Geuskens, T. Karnik, M. Khellah, J. Tschanz, and V. De, “Capacitive-coupling wordline boosting with self-induced vcc collapse for write vmin reduction in 22-nm 8t sram,” in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2012, pp. 234–236.
- [67] E. Karl, Y. Wang, Y.-G. Ng, Z. Guo, F. Hamzaoglu, M. Meterelliyoz, J. Keane, U. Bhattacharya, K. Zhang, K. Mistry, and M. Bohr, “A 4.6 ghz 162 mb sram design in 22 nm tri-gate cmos technology with integrated read and write assist circuitry,” *IEEE J. Solid-State Circuits*, vol. 48, no. 1, pp. 150–158, Jan 2013.
- [68] Y. Fujimura, O. Hirabayashi, T. Sasaki, A. Suzuki, A. Kawasumi, Y. Takeyama, K. Kushida, G. Fukano, A. Katayama, Y. Niki, and T. Yabe, “A configurable SRAM with constant-negative-level write buffer for low-voltage operation with $0.149\mu\text{m}^2$ cell in 32nm high-k metal-gate CMOS,” in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2010, pp. 348–349.
- [69] Y.-H. Chen, W.-M. Chan, W.-C. Wu, H.-J. Liao, K.-H. Pan, J.-J. Liaw, T.-H. Chung, Q. Li, G. Chang, C.-Y. Lin, M.-C. Chiang, S.-Y. Wu, S. Natarajan, and J. Chang, “A 16nm 128mb sram in high-k; metal-gate finfet technology with write-assist circuitry for low-vmin applications,” in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2014, pp. 238–239.
- [70] J. Chang, Y.-H. Chen, H. Cheng, W.-M. Chan, H.-J. Liao, Q. Li, S. Chang, S. Natarajan, R. Lee, P.-W. Wang, S.-S. Lin, C.-C. Wu, K.-L. Cheng, M. Cao, and G. Chang, “A 20nm 112mb sram in high-k; metal-gate with assist circuitry for low-leakage and low-vmin applications,” in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2013, pp. 316–317.
- [71] Y. Fujimura, O. Hirabayashi, T. Sasaki, A. Suzuki, A. Kawasumi, Y. Takeyama, K. Kushida, G. Fukano, A. Katayama, Y. Niki, and T. Yabe, “A configurable sram with constant-negative-level write buffer for low-voltage operation with $0.149/\mu\text{m}^2$ cell in 32nm high-k metal-gate cmos,” in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2010, pp. 348–349.
- [72] Y. Sinangil and A. Chandrakasan, “An embedded energy monitoring circuit for a 128kbit sram with body-biased sense-amplifiers,” in *IEEE Asian Solid State Circuits Conference*, Nov 2012, pp. 69–72.
- [73] Y. Sinangil, S. M. Neuman, M. E. Sinangil, N. Ickes, G. Bezerra, E. Lau, J. E. Miller, H. C. Hoffmann, S. Devadas, and A. P. Chandrakasan, “A Self-Aware Processor SoC using Energy Monitors Integrated into Power Converters for Self-Adaptation,” in *Symp. on VLSI Circuits (VLSI) Dig. Tech. Papers*, June. 2014, [Accepted].

- [74] A. Rahimi, L. Benini, and R. Gupta, "Analysis of instruction-level vulnerability to dynamic voltage and temperature variations," in *Design, Automation and Test In Europe (DATE) Dig. Tech. Papers*, March 2012, pp. 1102–1105.
- [75] Y. Ramadass, "An Energy Optimal Power Supply for Digital Circuits," Master's thesis, Massachusetts Institute of Technology, 2006.
- [76] Y. Ramadass and A. Chandrakasan, "Minimum energy tracking loop with embedded dc-dc converter enabling ultra-low-voltage operation down to 250 mv in 65 nm cmos," *IEEE J. Solid-State Circuits*, vol. 43, no. 1, pp. 256–265, Jan 2008.
- [77] "The MIT Angstrom Project: Universal Technologies For Exascale Computing." [Online]. Available: <http://projects.csail.mit.edu/angstrom/>.
- [78] J. Miller, H. Kasture, G. Kurian, C. Gruenwald, C. N. Beckmann, Celio, J. Eastep, and A. Agarwal, "Graphite: A distributed parallel simulator for multicores," in *High Performance Computer Architecture (HPCA), 2010 IEEE 16th International Symposium on*, Jan 2010, pp. 1–12.
- [79] L. Chang, D. Fried, J. Hergenrother, J. Sleight, R. Dennard, R. Montoye, L. Sekaric, S. McNab, A. Topol, C. Adams, K. Guarini, and W. Haensch, "Stable sram cell design for the 32 nm node and beyond," in *Symp. on VLSI Circuits (VLSI) Dig. Tech. Papers*, June 2005, pp. 128–129.
- [80] M. J. M. Pelgrom, A. C. J. Duinmaijer, and A. P. G. Welbers, "Matching Properties of MOS Transistors," *IEEE J. Solid-State Circuits*, vol. 24, no. 5, pp. 1433–1439, Oct. 1989.
- [81] Y. Sinangil, "Fault Tolerant, Low Voltage SRAM Design," Master's thesis, Massachusetts Institute of Technology, 2010.
- [82] "Predictive Technology Model (PTM)." [Online]. Available: ptm.asu.edu.
- [83] K. Zhang, K. Hose, V. De, and B. Senyk, "The scaling of data sensing schemes for high speed cache design in sub-0.18 μ m technologies," in *Symp. on VLSI Circuits (VLSI) Dig. Tech. Papers*, June 2000, pp. 226–227.
- [84] Y. Wang, U. Bhattacharya, F. Hamzaoglu, P. Kolar, Y. Ng, L. Wei, Y. Zhang, K. Zhang, and M. Bohr, "A 4.0 ghz 291mb voltage-scalable sram design in 32nm high-k metal-gate cmos with integrated power management," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2009, pp. 456–457,457a.
- [85] E. Karl, Y. Wang, Y.-G. Ng, Z. Guo, F. Hamzaoglu, U. Bhattacharya, K. Zhang, K. Mistry, and M. Bohr, "A 4.6ghz 162mb sram design in 22nm tri-gate cmos technology with integrated active vmin-enhancing assist circuitry," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2012, pp. 230–232.

- [86] T. Song, W. Rim, J. Jung, G. Yang, J. Park, S. Park, K.-H. Baek, S. Baek, S.-K. Oh, J. Jung, S. Kim, G. Kim, J. Kim, Y. Lee, K. S. Kim, S.-P. Sim, J. S. Yoon, and K.-M. Choi, "A 14nm finfet 128mb 6t sram with vmin-enhancement techniques for low-power applications," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2014, pp. 232–233.
- [87] A. Singh, M. Ciraula, D. Weiss, J. Wu, P. Bauser, P. de Champs, H. Daghighian, D. Fisch, P. Graber, and M. Bron, "A 2ns-read-latency 4mb embedded floating-body memory macro in 45nm soi technology," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2009, pp. 460–461,461a.
- [88] K. Nii, M. Yabuuchi, Y. Tsukamoto, Y. Hirano, T. Iwamatsu, and Y. Kihara, "A 0.5v 100mhz pd-soi sram with enhanced read stability and write margin by asymmetric mosfet and forward body bias," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2010, pp. 356–357.
- [89] H. Pilo, C. Adams, I. Arsovski, R. Houle, S. Lamphier, M. Lee, F. Pavlik, S. Sambatur, A. Seferagic, R. Wu, and M. Younus, "A 64mb sram in 22nm soi technology featuring fine-granularity power gating and low-energy power-supply-partition techniques for 37% leakage reduction," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb 2013, pp. 322–323.
- [90] P. Magarshack, P. Flatresse, and G. Cesana, "Utbb fd-soi: A process/design symbiosis for breakthrough energy-efficiency," in *Design, Automation and Test In Europe (DATE) Dig. Tech. Papers*, March 2013, pp. 952–957.
- [91] P. Flatresse, "Utbb fdsoi design & migration methodology," STMicroelectronics Technology R&D, Tech. Rep., Sept 2013.
- [92] N. Planes, O. Weber, V. Barral, S. Haendler, D. Noblet, D. Croain, M. Bocat, P. Sassoulas, X. Federspiel, A. Cros, A. Bajolet, E. Richard, B. Dumont, P. Perreau, D. Petit, D. Golanski, C. Fenouillet-Beranger, N. Guillot, M. Rafik, V. Huard, S. Puget, X. Montagner, M. A. Jaud, O. Rozeau, O. Saxod, F. Wacquand, F. Monsieur, D. Barge, L. Pinzelli, M. Mellier, F. Boeuf, F. Arnaud, and M. Haond, "28nm fdsoi technology platform for high-speed low-voltage digital applications," in *Symp. on VLSI Technology (VLSI) Dig. Tech. Papers*, June 2012, pp. 133–134.
- [93] O. Weber, F. Andrieu, J. Mazurier, M. Casse, X. Garros, C. Leroux, F. Martin, P. Perreau, C. Fenouillet-Beranger, S. Barnola, R. Gassilloud, C. Arvet, O. Thomas, J.-P. Noel, O. Rozeau, M. A. Jaud, T. Poiroux, D. Lafond, A. Toffoli, F. Allain, C. Tabone, L. Tosti, L. Brevard, P. Lehnen, U. Weber, P. K. Baumann, O. Boissiere, W. Schwarzenbach, K. Bourdelle, B.-Y. Nguyen, F. Boeuf, T. Skotnicki, and O. Faynot, "Work-function engineering in gate first technology for multi-vt dual-gate fdsoi cmos on utbox," in *Electron Devices Meeting (IEDM), 2010 IEEE International*, Dec 2010.

- [94] S. Narendra, A. Keshavarzi, B. Bloechel, S. Borkar, and V. De, "Forward body bias for microprocessors in 130-nm technology generation and beyond," *IEEE J. Solid-State Circuits*, vol. 38, no. 5, pp. 696–701, May 2003.
- [95] F. Arnaud, N. Planes, O. Weber, V. Barral, S. Haendler, P. Flatresse, and F. Nyer, "Switching energy efficiency optimization for advanced cpu thanks to utbb technology," in *Int. Electron Devices Meeting (IEDM) Dig. Tech. Papers*, Dec 2012.
- [96] C. Shin, M. H. Cho, Y. Tsukamoto, B.-y. Nguyen, B. Nikolic, and T.-J. K. Liu, "Sram yield enhancement with thin-box fd-soi," in *SOI Conference, 2009 IEEE International*, Oct 2009, pp. 1–2.
- [97] R. Tsuchiya, N. Sugii, T. Ishigaki, Y. Morita, H. Yoshimoto, K. Torii, and S. Kimura, "Low voltage (vdd $\tilde{0.6}$ v) sram operation achieved by reduced threshold voltage variability in sotb (silicon on thin box)," in *Symp. on VLSI Technology (VLSI) Dig. Tech. Papers*, June 2009, pp. 150–151.
- [98] A. Wang and A. Chandrakasan, "A 180mV FFT Processor Using Sub-threshold Circuit Techniques," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, 2004, pp. 292–293.
- [99] L. Svensson and J. Koller, "Driving a capacitive load without dissipating fcv^2 ," in *Low Power Electronics, 1994. Digest of Technical Papers., IEEE Symposium*, Oct 1994, pp. 100–101.
- [100] S. Younis and T. Knight, "Practical implementation of charge recovering asymptotically zero power cmos," in *Proceedings of the 1993 Symposium on Integrated Systems*, 1993, p. 234.
- [101] M. Khellah, D. Khalil, D. Somasekhar, Y. Ismail, T. Karnik, and V. De, "Effect of power supply noise on sram dynamic stability," in *Symp. on VLSI Circuits (VLSI) Dig. Tech. Papers*, June 2007, pp. 76–77.
- [102] R. Joshi, R. Kanj, S. Nassif, D. Plass, Y. Chan, and C.-T. Chuang, "Statistical exploration of the dual supply voltage space of a 65nm pd/soi cmos sram cell," in *Solid-State Device Research Conference, 2006. ESSDERC 2006. Proceeding of the 36th European*, Sept 2006, pp. 315–318.
- [103] A. Kawasumi and et al., "A Single-Power-Supply 0.7V 1GHz 45nm SRAM with An Asymmetrical Unit- β -ratio Memory Cell," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2008, pp. 382–383.
- [104] M. Sharifkhani and M. Sachdev, "Sram cell stability: A dynamic perspective," *IEEE J. Solid-State Circuits*, vol. 44, no. 2, pp. 609–619, Feb 2009.
- [105] M. E. Sinangil, "Low-Power and Application-Specific SRAM Design for Energy-Efficient Motion Estimation," Ph.D. dissertation, Massachusetts Institute of Technology, 2012.

- [106] T. Song, W. Rim, J. Jung, G. Yang, J. Park, S. Park, K. Baek, S. Baek, S. Oh, J. Jung, S. Kim, G. Kim, J. Kim, Y. Lee, K. Kim, S. Sim, J. Yoon, and K. Choi, "A 14nm FinFET 128Mb 6T SRAM with V_{MIN} -Enhancement Techniques for Low-Power Applications," *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, pp. 238–239, Feb. 2014.
- [107] E. Lau, J. E. Miller, I. Choi, Y. D., A. S., and A. A., "Multicore performance optimization using partner cores," in *HotPar*, 2011.
- [108] F. Kurdahi, A. Eltawil, Y.-H. Park, R. Kanj, and S. Nassif, "System-level sram yield enhancement," in *Int. Symp. on Quality Elec. Design (ISQED) Dig. Tech. Papers*, March 2006.
- [109] N. S. Kim, S. Draper, S.-T. Zhou, S. Katariya, H. Ghasemi, and T. Park, "Analyzing the impact of joint optimization of cell size, redundancy, and ecc on low-voltage sram array total area," vol. 20, no. 12, pp. 2333–2337, Dec 2012.
- [110] M. Yabuuchi, K. Nii, Y. Tsukamoto, S. Ohbayashi, Y. Nakase, and H. Shinohara, "A 45nm 0.6v cross-point 8t sram with negative biased read/write assist," in *VLSI Circuits, 2009 Symposium on*, June 2009, pp. 158–159.
- [111] "Leon3 Processor." [Online]. Available: <http://www.gaisler.com/index.php/products/processo>
- [112] "Leon3 Port for BEE2 and ASIC Implementation." [Online]. Available: http://cadlab.cs.ucla.edu/software_release/bee2leon3port/.
- [113] K. Marcinek, A. Luczyk, and W. Pleskacz, "Enhanced leon3 low power ip core for dsm technologies," in *Mixed Design of Integrated Circuits Systems, MIXDES*, June 2009, pp. 262–265.