# Characterization and Mitigation of Process Variation in Digital Circuits and Systems

by

Nigel Anthony Drego

B.S., Computer Engineering, University of California, Irvine (2001)
S.M., Electrical Engineering & Computer Science, Massachusetts
Institute of Technology (2003)

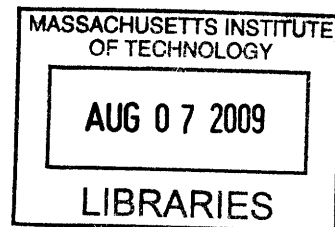Submitted to the Department of Electrical Engineering and Computer
Science
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2009

Author . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Department of Electrical Engineering and Computer Science
April 24, 2009

Certified by . . . . . . . . . . . . . . . . . . . . . . .
Duane Boning
Professor of Electrical Engineering & Computer Science
Thesis Supervisor

Certified by . . . . . . . . . . . . . . . . . . . . . . .
Anantha Chandrakasan
Professor of Electrical Engineering & Computer Science
Thesis Supervisor

Accepted by . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Terry P. Orlando
Chairman, Department Committee on Graduate Theses

# Characterization and Mitigation of Process Variation in Digital Circuits and Systems

by

## Nigel Anthony Drego

Submitted to the Department of Electrical Engineering and Computer Science
on April 24, 2009, in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy

## Abstract

Process variation threatens to negate a whole generation of scaling in advanced process technologies due to performance and power spreads of greater than 30-50%. Mitigating this impact requires a thorough understanding of the variation sources, magnitudes and spatial components at the device, circuit and architectural levels. This thesis explores the impacts of variation at each of these levels and evaluates techniques to alleviate them in the context of digital circuits and systems.

At the device level, we propose isolation and measurement of variation in the intrinsic threshold voltage of a MOSFET using sub-threshold leakage currents. Analysis of the measured data, from a test-chip implemented on a $0.18\mu m$ CMOS process, indicates that variation in MOSFET threshold voltage is a truly random process dependent only on device dimensions. Further decomposition of the observed variation reveals no systematic within-die variation components nor any spatial correlation.

A second test-chip capable of characterizing spatial variation in digital circuits is developed and implemented in a $90nm$ triple-well CMOS process. Measured variation results show that the within-die component of variation is small at high voltages but is an increasing fraction of the total variation as power-supply voltage decreases. Once again, the data shows no evidence of within-die spatial correlation and only weak systematic components. Evaluation of adaptive body-biasing and voltage scaling as variation mitigation techniques proves voltage scaling is more effective in performance modification with reduced impact to idle power compared to body-biasing.

Finally, the addition of power-supply voltages in a massively parallel multicore processor is explored to reduce the energy required to cope with process variation. An analytic optimization framework is developed and analyzed; using a custom simulation methodology, total energy of a hypothetical 1K-core processor based on the RAW core is reduced by 6-16% with the addition of only a single voltage. Analysis of yield versus required energy demonstrates that a combination of disabling poor-performing cores and additional power-supply voltages results in an optimal trade-off between performance and energy.

Thesis Supervisor: Duane Boning
Title: Professor of Electrical Engineering & Computer Science

Thesis Supervisor: Anantha Chandrakasan
Title: Professor of Electrical Engineering & Computer Science

# Acknowledgments

Though a PhD is primarily the undertaking of a single person, it cannot be done without the aid of many others. I have been infinitely blessed to be surrounded by family, friends, labmates, colleagues and professors who have provided support and aid throughout this PhD.

I would first like to thank my research advisors, Prof. Duane Boning and Prof. Anantha Chandrakasan who have provided not only technical and research guidance, but practical advice on many other topics as well. I couldn't ask for better advisors over the course of the past four years — I have learned a lot from both of you and I only hope that throughout my career, I am able to pass along a fraction of the wisdom and learnings you have passed on to me.

Two other professors provided excellent support and guidance over the last year of my PhD. The third member of my thesis committee, Prof. Anant Agarwal, pointed us in the right direction when it came to understanding the future of multicore processors and the critical research needs to enable continued performance scaling. Prof. Devavrat Shah's masterful command of mathematics and optimization thoroughly aided me during the last part of my PhD. I am grateful for Prof. Shah's time and patience in helping me to understand optimization and mathematical bounding of the optimization problem I faced.

My family has been a continuing source of love and support throughout my life, no less so during my PhD. To Mom, Dad, Roulla, Priya, Pathy and Naveen, I know how proud you are of me and the culmination of this work and I can't thank you enough for the unconditional backing you have always provided and continue to provide. To Vid, arguably you deserve this PhD even more than I do — it has been you that has been my source of sanity. You are my rock of support both in times of difficulty and otherwise, and a counterpoint, providing a cool head in times of frustration and a "kick in the rear" when I've lacked motivation. Without you, it's impossible to imagine this ever being undertaken or completed — Thank You!

The value of a good set of friends should never be underestimated. Over the past

four years, I have been extremely fortunate to meet and be in the company of a stellar group of friends. Anand and Nammi, you are our extended family here. You took us in when we were searching for housing out here and since then you've provided both Vid and myself with so many laughs, good times, diversions and some amazing chocolate chip cookies! Naveen and Yogesh, you both are exceptional friends who have helped me so much, both technically and otherwise. Coffee hour has been such a critical part of my PhD and I do not look forward to the day that ends. There are many other friends who must be acknowledged for their many contributions to my technical understanding and/or well-being over the past four years, including Karthik Balakrishnan, Daihyun Lim, Karen and Andy Gettings and Anita Misri as well as Nicholas Velastegui and Steve Pfeiffer from my days at Intel. In addition, I can't thank Debb Hodges-Pabon enough for being another mother and such a wonderful confidant and human-being. Debb, you truly do make everyone's MTL experience so much better.

I am grateful to all the past and present members of both the Boning and Anantha research groups, with whom I've shared many an interesting conversation and who have helped me in many capacities.

# Contents

# List of Figures

13

# List of Tables

# Chapter 1

# Introduction

The continued integration and compression of the modern electronics we often take for granted has been the result of continuous transistor scaling. A device like the iPhone, that combines communication, entertainment, navigation and personal information management, is simply impossible to construct in its given form factor without integrating hundreds of millions, if not billions, of transistors, each with dimensions of $90nm$ or smaller. Modern microprocessors are comprised of transistors with electrical properties that can be described by numbers of atoms present or thicknesses counted by the number of atomic layers stacked. Perhaps more astonishing is that this steadfast increase in transistor density, known as Moore's Law [14], has provided faster and more powerful electronics with constant or even decreasing cost.

In 1965, Gordon Moore observed that the density of transistors on a die increased by a factor of two every 18 months - an observation that was quickly dubbed "Moore's Law" [14]. The decreased cost and increased performance associated with increased density makes an electronics consumer believe that a particular device purchased today will either be cheaper or have more features for the same cost in the future (i.e., the consumer implicitly takes Moore's Law to be a law rather than an observation). Upon closer inspection however, there have been and continue to be many a technological hurdle to overcome in advancing Moore's *observation*. One of the most worrisome challenges in the current decananometer era of semiconductors is variation: two nominally identical transistors, when fabricated, will vary in many respects

(Figure 1-1).



Figure 1-1: Two "nominally" identical transistors that are physically different due to a variety of causes.

This process variation is not new to manufacturing lines: typically known as process tolerance[1], it is an established concept in a wide range of manufacturing processes, from biological to mechanical and agricultural to electrical, and including semiconductor manufacturing. However, in most cases the magnitude of the variation is small relative to the nominal design parameters — with appropriate process control, these variations do not significantly impact the design nor the operation of the manufactured product. Until recently, this was the case in semiconductors as well: product-impacting variations were primarily due to and dominated by yield-loss defects and were mitigated or eliminated predominately with improved process control. Lately, however, this situation has deteriorated rapidly due to increasingly limited controllability of individual process modules operating at the limits, and in some cases beyond the originally intended limits.

To illustrate the impact of variation on actual products, Figure 1-2 plots the normalized distributions of frequency and standby leakage of Intel microprocessors on a single wafer. Parameter variations result in greater than 30% frequency spread and 20X variation in chip leakage. The large frequency spread necessitates expensive frequency binning in which each chip is tested to determine its maximum frequency and power before it can be sold — an often expensive, time-consuming process. Moreover,

---

[1]Process tolerance is the allowable variation in the parameters of a manufactured item that does not adversely affect the stated performance of that item.

as the standby leakage component of power increases as a fraction of the total power, 20X variation in leakage currents can mean that even if leakage is nominally only 1% of the total power, with variation it can be as much as 20%. In reality, variation in leakage power can result in variation of total power by as much as 50% [15].



Figure 1-2: Frequency and leakage variations of Intel microprocessors on a single wafer [1]

As a result, yield is affected by parameter variations: chips that operate too slowly with high standby leakage power, or those that have high performance but are above the power envelope, must be discarded. Though microprocessors often represent extreme examples of semiconductor engineering, the problem is more generally valid: performance and power are significantly impacted by unmitigated parameter variation resulting in parametric yield loss. This poses a challenge that requires careful analysis and a paradigm shift as device engineers, circuit designers and system architects all must now consider process variation during the course of technology and product design and development.

# 1.1   Thesis Organization

Effective mitigation of process variation requires both understanding and characterization of its effects at all levels of design. As will be seen in Section 2.5, considerable work has been carried out in many areas of the variation spectrum. Work in one important area, the characterization of spatial variation and its implications on digital circuits and systems, has been trailing, with only piecemeal contributions in very specific contexts. This thesis provides a comprehensive, bottom-up analysis of within-die spatial process variation in the context of digital circuits and systems.

We begin in Chapter 2 by putting variation in context. A summary of variation sources in key process steps provides the basis for evaluating the impact of variation on individual transistor parameters. To better characterize and model variation, we explain how to decompose variation, particularly at the within-die level, drawing a distinction between systematic and spatially (un)correlated components. This allows evaluation of the impact this decomposition has to the modeling and design of digital circuits and systems. Background and related work are also provided to conclude the chapter.

Characterization and analysis of variation at the device level is undertaken in Chapter 3, by specifically focusing on the intrinsic threshold voltage, $V_{T_o}$, of a transistor. Circuit designers are most often concerned about variation in channel length and $V_T$. With comprehensive spatial analysis of channel length variation provided by Friedberg et al. in [16], similar analysis of $V_T$ variation became crucial. A test-chip capable of characterizing spatial variation of $V_{T_o}$ by measuring the sub-threshold currents of thousands of transistors is implemented and measured. Analysis of the data is performed to characterize both within-die spatial variation as well as any systematic spatial patterns repeatable from die-to-die. In combination with knowledge of spatial variation in channel length, Monte Carlo simulations are performed to understand the implications at the circuit level.

Chapter 4 focuses on abstracting away device parameter variation and understanding spatial variation of common digital circuits. Here we detail a test-chip

architecture consisting of replicated blocks of digital circuits and high-precision measurement circuitry, implemented in a $90nm$ CMOS technology, that facilitates such understanding. The measured data from the test-chips are analyzed and compared to the Monte Carlo simulation results in Chapter 3. The results of Adaptive Body Biasing (ABB) as a potential mitigation scheme are also discussed and show that in advanced technologies Adaptive Voltage Scaling (AVS) is more effective.

Using the results and insights from the test-chips in Chapters 3 and 4, mitigation techniques at the architecture level are explored in Chapter 5. Specifically, we focus on future multi-core processors where the potential impact of within-die variation on both performance and power is significant. An analytic, mathematical approach is taken to choosing optimal power-supply voltage values in order to reduce energy. Using such an approach enables substantial energy savings, compared to a "worst-case" approach of using a single high-valued voltage, while meeting both yield and performance constraints, including core performance homogeneity, that are crucial to system designers and operating system architects.

Finally, Chapter 6 concludes with a high-level summary of this thesis, as well as suggestions for future research as informed by the findings outlined in the previous chapters.

## 1.2 Thesis Contributions

The major findings and contributions of this thesis are:

- Proposal of a measurement technique utilizing sub-threshold currents to effectively isolate and extract variation in the intrinsic threshold voltage of MOSFET transistors.

- Design and implementation of a test-chip using the proposed measurement technique to characterize variation in an array of $\sim 140K$ transistors.

- Design and implementation of an architecture to quantify spatial variation in digital circuits, including a high-speed, arbitrarily fine-resolution delay mea-

surement technique based on random-sampling and using only common digital components.

- Measurement analysis and decomposition of within-die variation for both MOS-FET threshold voltage and digital circuit performance into systematic (position-dependent), spatially correlated (distance-dependent) and spatially uncorrelated components. In both cases, within-die variation is primarily random and spatially *un*-correlated (i.e., no spatially correlated component). In the case of digital circuit performance, a systematic component is identifiable but is small in magnitude relative to other components.

- Analysis (both through simulation and measured data) showing that variation sensing techniques included in mitigation schemes are highly dependent on the performance mode of a circuit, posing unique challenges for systems that scale between high-performance and low-power operating modes.

- Proposal, implementation and analysis of an analytic framework and algorithm for robust optimization and minimization of "variation-induced energy overhead" in massively parallel multi-core processors.

- A custom simulation methodology enabling the above framework and algorithm to be employed on a real core. These simulations demonstrated that a single additional power-supply voltage can reduce total energy consumption considerably in the face of variation.

# Chapter 2

# Process Variation In Context

Process variation is increasingly becoming a limiting factor in both IC design and manufacturing [17], as nearly every step in the IC manufacturing process introduces variation in the end device. This chapter explores the sources of semiconductor process variation, their impact, and related work in the field. We begin in Section 2.1 by enumerating the most significant of these variation sources, and discuss their impact on individual transistor parameters in Section 2.2. This is followed in Section 2.3 by a discussion of their spatial dependencies, and by analysis of impact to modeling and design of circuits and systems in Section 2.4. Finally, Section 2.5 summarizes previous research on variation in semiconductor manufacturing.

## 2.1   Sources of Variation



Figure 2-1: Cross-section of a MOSFET

Figure 2-1 depicts the lateral cross-section of a MOSFET in its most simple, ideal form and Figure 2-2 shows a cross-sectional view of an entire integrated circuit and the major steps involved in the fabrication of such a circuit. In Deep-Sub-Micron (DSM) CMOS, each of these requires one or more unit process steps. For example, formation of one of the two twin-well implants involves depositing or thermally growing an oxide layer, spinning on photoresist, lithographically patterning the photoresist to define the well area, developing away the exposed photoresist over the defined well areas, implanting the appropriate dopant species and then removing the resist and oxide layer. This entire procedure is then repeated for the other well. Variations of this procedure are used for each of the first thirteen steps listed in Figure 2-2. When forming the Shallow Trench Isolation (STI) and interconnect layers, additional polishing steps are required to ensure a smooth, uniform surface on which subsequent layers can be fabricated. Modern product designs require 100+ individual process steps to fully fabricate the entire CMOS stack.

1. Twin-well Implants
2. Shallow Trench Isolation
3. Gate Structure
4. Lightly Doped Drain Implants
5. Sidewall Spacer
6. Source/Drain Implants
7. Contact Formation
8. Local Interconnect
9. Interlayer Dielectric to Via-1
10. First Metal Layer
11. Second ILD to Via-2
12. Second Metal Layer to Via-3
13. Metal-3 to Pad Etch
14. Parametric Testing



Figure 2-2: Cross-sectional view of a CMOS integrated circuit with major steps needed for fabrication [2].

A number of these process steps can be highlighted as major sources of variation [18]: 1) sub-wavelength lithography, 2) plasma etch, 3) ion implantation and annealing, and 4) chemical-mechanical polishing (CMP). Depending on the feature being fabricated, each process step affects subsequent transistor and interconnect parameters in differing manners: variations in lithography, etch and CMP affect the physical dimensions of transistors and the wires and vias that constitute the interconnect between transistors. However, ion implantation and annealing directly affect the molecular make-up of transistors. Furthermore, the significance and impact of variation in a particular process step is highly dependent on not only the feature being fabricated, but also the application in which the fabricated transistor is used. For example, variation in the size of the source/drain area of a transistor may impact the overall performance far greater when that transistor is used in an analog versus digital application.

As critical dimensions continue to decrease, process variations become increasingly worrisome due to decreasing depth-of-focus of sub-wavelength lithography, line-edge roughness, random discrete dopant fluctuation, stress effects, and oxide thickness ($t_{ox}$) fluctuation. In the following sub-sections, we identify some of the more significant sources of variation in each of the four process steps outlined above.

## 2.1.1 Lithography

Lithography is the process of exposing a light-sensitive material (photoresist, or just "resist") to define the critical physical dimensions of a semiconductor structure. Until the $180nm$ node, the wavelength of light used to pattern these critical dimensions scaled with the smallest of the dimensions to be patterned. In this regime lithography-induced variations were a result of lens imperfections, mask errors, illumination non-uniformity, and contributions arising from resist non-uniformities [19]. At the $180nm$ node, scaling of the wavelength of light used for patterning ceased at $193nm$ due to increased cost of lithography technology, materials, and equipment development and deployment. The resulting lithographic defocus causes both systematic and random line-width variations [20], mitigated to some degree by so-called Resolu-

tion Enhancement Techniques (RETs) such as Optimal Proximity Correction (OPC), Sub-Resolution Assist Features (SRAF) and phase-shifted mask lithography [21].

Perhaps more alarming are random variations in line edges, known as Line-Edge Roughness (LER) and depicted in Figure 2-3, resulting in local variations in line-width [22]. The sources of this variation are still under discussion, but conjectures include shot noise of the energy of the light illuminating the resist, solubility and size of resist polymer particles, and local variations in the chemistries that make up chemically amplified resists [23].



(a) Nominal transistor  (b) Transistor with LER and RDF  (c) High-resolution microscopy of LER  (d) LER affecting interconnect

Figure 2-3: Line-edge Roughness (LER) [3]

Immersion lithography, extreme ultra-violet (EUV) lithography and improved resist materials may aid in improved control of physical gate and line dimensions. However, to date, only immersion lithography is commercially viable despite years of research related to EUV and improved resists.

## 2.1.2  Plasma Etch

After lithographic patterning, plasma etching is used to etch away unneeded areas of polysilicon, in the case of transistor gates, or an insulator, such as silicon dioxide or a low-k dielectric, in the case of interconnect. Variations in process conditions such as chamber temperature and pressure, RF power, electrode spacing, and gas flows often result in variations if not properly controlled using statistical process control [24]. Product layout, local chemistry non-uniformities and process conditions,

and lithography-induced variations also result in etch non-uniformities, giving rise to variation in side-wall profiles (varying slope of side-walls), line-widths, and thicknesses (due to variation in etch depth). More detailed characterization and analysis of variations due to plasma etch can be found in the work of Abrokwah [25].

## 2.1.3 Ion Implantation & Annealing

Creation of transistors involves doping them with ions to define the type of transistor (PMOS or NMOS). The substrate, as well as other components of the transistor such as the highly-doped source/drain regions, are doped with different ion species (e.g. B, As, P). These ions are accelerated at high energy into the wafer during the ion implantation step and are then "activated" by heating the wafer (annealing) in order to ensure the implanted ions are properly substituted within the existing crystal structure of the underlying silicon substrate.

Once again, local and global process conditions such as implant energy and dose, tilt angle and temperature profiles all result in variations of the implanted ions. Layout features and proximity effects, such as distance to well edge, can also affect uniformity of ion implantation [18]. In the most advanced process technologies (e.g. $45/65nm$ technology nodes), device volumes are so small that only several tens to low hundreds of dopant atoms are needed within the channel area, directly underneath the gate, for the required doping concentrations. Due to the small numbers, variation in the dopant counts and even the placement of the atoms within the transistor body is of significant concern, as regions of a single transistor will experience different local doping concentrations. This variation mechanism is known as Random Dopant Fluctuation (RDF) and was brought to light as early as 1975 by Keyes [26]. The right of Figure 2-4 depicts RDF where the black dots in the channel are countable dopant atoms [4]. It is easy to see that by changing the number or even the placement of the atoms the electrical performance of the transistor can be greatly impacted. As transistor volumes continue to shrink, without individual placement of dopant atoms, RDF is unavoidable due to the decreasing absolute number of dopants required.

As RDF and other variation mechanisms arising from ion implantation and an-

Figure 2-4: Intel simulation of Random Dopant Fluctuation (RDF) [4].

nealing become increasingly significant in modern processes, many have suggested moving to significantly different transistor structures, such as fully depleted devices (e.g., ultra-thin body or FinFET devices) [10]. However, acceptance of and transition to radically different device structures is both technologically and economically difficult given the dominance and proven abilities of lateral MOSFETs. Consequently, the most common approach to mitigating this type of variation is to increase device size which reduces relative variation as the number of dopant atoms necessarily increases with device size. Such a solution is fundamentally incompatible with further transistor scaling, making it unsustainable in the long term and requiring solutions either at the process or design levels (i.e., improved process modules for the doping step, new device structures not requiring doping and/or circuit design that is robust to device variation).

### 2.1.4 Chemical-Mechanical Polishing (CMP)

CMP is used to achieve smooth and planar surfaces from which subsequent layers are able to be fabricated. Decreasing depth-of-focus in modern lithography systems underscores the need for exquisite planarity and without such planarity, features to be patterned may be out of focus due to surface height fluctuations (nanotopography) [27]. This results in subsequent lithographic variations as described in Section 2.1.1. However, CMP is not a variation-free process itself: it is a significant

source of systematic variation resulting from both process conditions, including variations in down force, rotational speed, pad conditioning, and temperature as well as designed feature sizes and pattern dependencies [28]. The primary effects of CMP variation are shown in Figure 2-5, where copper lines can be "dished" and inter-layer dielectrics eroded, causing variation in copper line thicknesses.



Figure 2-5: Dishing of copper and erosion of inter-layer dielectrics in CMP.

Mitigation strategies for variation resulting from the CMP process module began with improved process control by using feedback from the process itself to guide when the polishing should end [29]. More recently, mitigation strategies have focused on improved modeling and design modification: since variations arising from the CMP process tend to be limited to pattern dependencies and features sizes, appropriate modeling of the process and product design can reveal areas of particular susceptibility to the types of variation depicted in Figure 2-5. With this information, automated mitigation strategies have been devised to enforce or adjust pattern densities (e.g., by using design rule constraints or automated "dummy fill" insertion) to dramatically reduce a design's susceptibility to CMP-caused variation [30].

## 2.1.5 Other Variation Sources

The four process modules described above contribute significantly to overall process variation, but variation is by no means limited to these four modules. Other process steps that are sources of variation include gate oxidation, polysilicon and nitride deposition, and metallization, all of which can result in variation in film thicknesses,

with varying degrees of impact to the fabricated transistor. Film thickness variation in polysilicon and interconnect metal are typically mitigated using CMP. However, as the previous section described, this can be a source of variation as well.

Variations in the gate oxidation step are primarily wafer-to-wafer due to differences in chamber temperature and length of time in the chamber. Within-wafer variations, due to temperature-induced stresses, lamp configuration and convective cooling, are corrected with better tool design as well as improved process control [31]. As a result, gate oxide and nitride layers are generally well controlled at the process module level but as dimensions continue to reduce, even small variations are amplified.

While we have described effective solutions for many of the variation sources described above, solutions for improved control of random discrete dopant fluctuation or oxide thickness at the manufacturing level do not exist, meaning these are issues circuit designers and system architects must increasingly be aware of and learn to deal with. Dealing with variation requires understanding how such fluctuations affect device properties, circuit performance and their architectural implications.

## 2.2   Impact on Transistor Parameters

Each of the variation sources highlighted above (and others) impacts the electrical properties of transistors and interconnect in unique and often subtle ways. These effects are best understood in the context of transistor performance. In a typical digital integrated circuit, a transistor either charges or discharges a capacitive load, and the time required to do so determines the performance of the transistor. This time is a function of the capacitance being driven, the voltage to which it must be driven and the current used to drive it, as shown in Eq. 2.1. For simplicity, we use the ideal I-V equation for a MOSFET in the saturation regime as shown in Eq. 2.2, where $\mu$ is the mobility of a charge carrier through the device, $C_{ox}$ is the gate oxide capacitance, $W$ and $L$ are respectively the width and length of the transistor, $V_T$ is the device threshold voltage and $V_{GS}$ is the bias between gate and source. Though this equation is idealized and neglects important details in modern transistors, it is

sufficient to illustrate the impacts that the variation sources mentioned above have on a transistor.

$$t_d = \frac{C_{load}V_{DD}}{I} \tag{2.1}$$

$$I_D = \frac{1}{2}\mu C_{ox}\frac{W}{L}\left(V_{GS} - V_T\right)^\alpha \tag{2.2}$$

$$t_d = \frac{C_{load}V_{DD}}{\frac{1}{2}\mu C_{ox}\frac{W}{L}\left(V_{GS} - V_T\right)^\alpha} \tag{2.3}$$

Table 2.1 shows the MOSFET parameters and relevant process modules that directly affect each of those parameters. It is clear that a single process module can affect multiple transistor parameters, and thus decoupling the effects of one variation source from another are difficult. Nevertheless, we now explore variation from the perspective of the device, in particular each of the parameters listed in the table.

| MOSFET Parameter | Relevant Process Module(s) |
|---|---|
| $\mu$ | Ion implantation, annealing, diffusion, nitride deposition |
| $C_{ox}$ | Gate oxidation |
| $W, L$ | Lithography, etch |
| $V_T$ | Ion implantation, annealing, gate oxidation, (lithography, etch) |

Table 2.1: Process modules affecting various transistor parameters.

## 2.2.1 Mobility ($\mu$)

Mobility refers to the ease which charge carriers (electrons or holes) can travel through the channel of a MOSFET in response to an applied electric field. It is mathematically defined as in Eq. 2.4, where $q$ is the electronic charge, $\tau_c$ is the mean free time between carrier collisions, and $m_{n,p}$ is the effective mass of either an electron (n) or hole (p). However, in practice, mobility is given as a function of the doping concentration as shown in Figure 2-6, since the doping concentration determines the mean free time between collisions, and to a lesser degree, the effective mass. In modern processes, stress engineering in the form of nitride liners and silicon germanium source/drains, also affects mobility by either stretching or compressing the silicon lattice to decrease

31

the effective mass of a particular charge carrier [32].

$$\mu_{n,p} = \frac{q\tau_c}{2m_{n,p}}$$

(2.4)

Any process step which affects doping concentration or stress will necessarily affect transistor mobility. Therefore, ion implantation and annealing directly affect mobility as these process steps primarily determine doping concentrations. However, as seen in Figure 2-6, since doping concentration is on a log scale and typically does not vary by orders of magnitude from one transistor to another, the impact that ion implantation, annealing and other process modules that determine doping concentration have on mobility is relatively small.



Figure 2-6: Mobility as function of doping in Si [5].

Intentional and unintentional stresses, whether by stress engineering or proximity to STI, can have large impacts on transistor mobility. Mobility improvements greater than 10% over unstrained silicon have been reported as strain engineering has matured [33]. Even unintentional stresses due to STI proximity can cause within-die mobility variations on the order of a few percent depending on transistor distance to the STI edge [34]. Recent characterization of mobility in advanced processes indicates relatively large variation, 21% $\frac{\sigma}{\mu}$, and may be due to fluctuations in the intentional stresses introduced in these processes [35].

## 2.2.2 Oxide Capacitance ($C_{ox}$)

Gate oxide capacitance is the capacitance between the gate stack (polysilicon and silicon dioxide) and the inverted channel of the MOSFET. Eq. 2.5 shows that the oxide capacitance is a function only of the oxide thickness ($t_{ox}$) and the dielectric constant of silicon dioxide or other gate insulator.

$$C_{ox} = \frac{\epsilon_{ox}}{t_{ox}} \tag{2.5}$$

As mentioned in Section 2.1.5, gate oxidation (thermal growth of silicon dioxide or silicon nitride) is a relatively well controlled process step. However, with gate oxide thicknesses on the order of five atomic layers ($\approx 1nm$), even small variations of one atomic layer have the potential to greatly impact not only oxide capacitance, but also threshold voltage and mobility [13]. With $SiO_2$ gate oxide, variations of a single monolayer (approximately $0.2nm$) are typical and result in 20% shifts in oxide thickness. Furthermore, though not shown in Eq. 2.2, oxide thickness exponentially impacts gate currents due to Fowler-Nordheim tunneling [36]. As a result, variation in oxide thickness can have significant and pernicious effects on idle leakage power in modern devices. To improve this situation, Intel has recently begun using a "high-K" gate dielectric, hafnium dioxide ($HfO_2$), to allow for thicker gate oxides while maintaining oxide capacitance and gate control over the channel, but reducing gate leakage by three orders of magnitude [37, 38]. Moving to a new gate oxide material not only reduces gate leakage currents but reduces the impact of variability on $C_{ox}$ due to the much larger physical oxide thickness. Nevertheless, variations in the interfacial oxide of "high-K" stacks do still occur and can affect performance [39].

## 2.2.3 Transistor Dimensions ($W$, $L$)

The previous section highlighted the importance of one of the smallest dimensions of a transistor, the gate oxide thickness. Eq. 2.2 shows that the width ($W$) and length

($L$) of a transistor are critical in determining the current through that transistor.[1] $W$ must be increased or $L$ decreased to increase current and thus performance. Since decreasing $L$ also reduces load capacitances and increases transistor density, scaling has continually reduced $L$ in the pursuit of increased performance, making it the most critical dimension in a transistor today.

Lithographic patterning and etching define both of these dimensions. However, since $W$ is always larger than $L$, only variation in channel length is typically of concern (with the exception of the smallest width devices). Eq. 2.3 shows that the delay of a transistor is directly proportional to the channel length, so any variation in channel length will be directly reflected in transistor delay. As transistor lengths have decreased well below the wavelength of light patterning them, relative variation in channel length has increased. This is evident in International Technology Roadmap for Semiconductors (ITRS) projections shown in Figure 2-7, where the 2001 and 2003 projections were 10% into the forseeable future, but in 2005 and 2007 the $\frac{3\sigma}{\mu}$ projections increased to 12% which, without mitigation, portends a corresponding increase in performance variability. It would not be surprising to see another increase in the 2009 edition of the roadmap. Although not mentioned above, variation in



Figure 2-7: ITRS projections for channel length variation [6, 7, 8, 9].

---

[1]Not shown by Eq. 2.2 are second-order dependencies on channel length that affect threshold voltage ($V_T$) and device leakage currents but which will be discussed to some degree in Section 2.2.4 and Chapter 3.

channel length also results in threshold voltage variations due to the Drain-Induced Barrier Lowering (DIBL) phenomenon, which has been known for over 25 years but has had a more dramatic effect as channel lengths scaled below $100nm$ [40]. Modern processes typically show measured channel length variation of 3-4% $\frac{\sigma}{\mu}$, consistent with the ITRS projections [35, 1].

## 2.2.4   Threshold Voltage ($V_T$)

The threshold voltage of a MOSFET is the gate-to-source bias ($V_{GS}$) that results in a channel forming just under the gate, allowing current conduction from source to drain of the transistor. In an ideal long-channel MOSFET, the threshold voltage is determined by only the doping concentration ($N_{ch}$) and the oxide capacitance ($C_{ox}$) as shown in Eq. 2.6, where $V_{FB}$, the flatband voltage, and $\phi_{F_p}$, the Fermi Potential of the substrate, are dependent only on the doping concentration, and $\gamma$ is dependent on both the doping concentration and oxide capacitance. In short-channel devices, effects such as DIBL result in $V_T$ being additionally dependent on the channel length ($L$), source/drain junction depths ($x_j$), and stresses [5], meaning that a large fraction of process steps can potentially affect the value of $V_T$.

$$V_{T_o} = V_{FB} + 2\phi_{F_p} + \gamma\sqrt{2\phi_{F_p} + V_{SB}} \qquad (2.6)$$

Owing to this "susceptibility" and the intrinsic random variability of RDF, $V_T$ is one of the least controlled transistor parameters, with $3\sigma$ variations on the order of 30% or more of mean [13]. However, it is also one of the most studied. Pelgrom et al. showed that variation in $V_T$ is a function of the area of the device as in Eq. 2.7 [41]. As scaling has led to smaller and smaller device dimensions, control of $V_T$ has gotten progressively worse, as evidenced in Figure 2-8. Since the mean value of $V_T$ has been reduced over those technology generations, relative variation ($\frac{\sigma}{\mu}$) increases even more rapidly. For a minimum size device in a high-performance $45nm$ process where the

$\overline{V_T} = 0.25V$, we find a large $\frac{\sigma}{\mu}$ of approximately 20%.

$$\sigma^2_{V_T} = \frac{A^2_{V_T}}{L_{eff}W_{eff}} + S^2_{V_T}D^2 \tag{2.7}$$

More recently, Asenov et al. refined Pelgrom's model and formulated the empirical



Figure 2-8: $\sigma V_T$ has increased with continued technology scaling [10].

expression shown in Eq. 2.8 which describes the standard deviation of $V_T$ in relation to the fundamental parameters of doping concentration, oxide thickness and transistor dimensions [42].

$$\sigma_{V_T} = 3.19 \times 10^{-8} \frac{t_{ox}N^{0.4}_{ch}}{\sqrt{L_{eff}W_{eff}}} \tag{2.8}$$

Looking at Eq. 2.3, the impact that variation in $V_T$ has on performance can be substantial due to the quadratic dependence on $V_T$. In many cases the gate-overdrive, $V_{GS} - V_T$, is large enough that $V_T$ variation can be mitigated, but in low-power designs where $V_{DD}$, and thus $V_{GS}$, are not much greater than $V_T$, the impact is large. Moreover, as will be shown in Chapter 3, device leakage currents are exponentially dependent on $V_T$, so $V_T$ variation has considerable impact on idle power.

## 2.2.5  Device Impact Summary

The impact of variation in each of the above transistor parameters differs depending on a variety of factors, including circuit implementation, logic style and region of

operation. Furthermore, the impact "signature" of variation in a particular parameter is different for different metrics. For the reader's ease, tables summarizing the impact of variation in the above transistor parameters are provided based on data found in [13] from Monte Carlo simulations of common digital blocks, such as adders, implemented in a $90nm$ process. Table 2.2 shows the variability in delay and power as a function of different circuit styles, while Table 2.3 breaks down the overall variation into contributions from individual parameters. Variation in $V_T$ and channel length contribute most heavily to overall variation; this will be a recurring message throughout this thesis, especially in Chapters 3 and 4. However, it should also be noted that temporal sources of variation, especially $V_{DD}$ fluctuations, which are not included in this summary also significantly contribute to the overall variation.

| Circuit style | Delay Variability $\left(\frac{\sigma}{\mu}\right)$ (%) | Power Variability $\left(\frac{\sigma}{\mu}\right)$ (%) |
|---|---|---|
| Static CMOS | 6.1% | 4.1% |
| Pulsed-Static CMOS | 6.5% | 5% |
| Domino | 6.6% | 4.3% |

Table 2.2: Variability in delay and power based on circuit-style [13].

| Parameter | Delay $\left(\frac{\sigma}{\mu}\right)$ (%) | Power $\left(\frac{\sigma}{\mu}\right)$ (%) |
|---|---|---|
| $t_{ox}$ | 1-2% | 1-2% |
| W | $\ll 1\%$ | 0.5-1% |
| L | $\sim 3\%$ | $< 2\%$ |
| $V_T$ | 2.5-6% | 1.75-4.75% |

Table 2.3: Contributions of variation from individual parameters over various circuits and circuit styles, summarized from data in [13].

## 2.3 Decomposing Variation

The physics of each variation source results in both temporal and spatial dependencies of the resulting variation: the statistics of the variation are dependent on time of fabrication (e.g., fabs continuously tweak process parameters to improve yield) as well as distance and length scales (within-die variation versus die-to-die or wafer-to-wafer). Temporal dependencies can also include environmental variation sources such

Figure 2-9: Decomposition of process variation into different length scales.

as power-supply variations and noise, or longer-term degradation/reliability mechanisms like Negative-Bias Temperature Instability (NBTI) which causes device threshold voltages to shift over the lifetime of a product [43]. However, this thesis will intentionally omit discussion of temporal decomposition of variation and instead focus on spatial decomposition.

Spatial decomposition of variation is natural in the context of products smaller than the wafers they are fabricated on; process engineers and circuit designers often want to know how variation statistics change at different length scales. For example, process engineers are often concerned with wafer and die yields (number of sellable die on good wafers) and may thus be interested in wafer-to-wafer variations or even within-wafer trends. On the other hand, circuit designers may be more concerned with how two transistors (or circuits) close to each other on a die vary with respect to each other (within-die variation), versus how two transistors (or circuits) vary from one die to another (die-to-die variation). Given the varying interests and concerns, spatial decomposition is necessary for and critical to a complete understanding of variation.

Variation of a parameter, $P$, is typically decomposed into lot-to-lot (L2L), wafer-

to-wafer (W2W), die-to-die (D2D), within-die (WID), and random (RND) components as depicted in Figure 2-9. Mathematically, this can be stated as in Eq. 2.9 where the final value of $P$ is the sum of the components above. $\mu_{X,Y}$ represents any variation that results from a circuit being placed at a particular location on the die, typically referred to as within-die systematic components of variation, while $\epsilon$ represents the residue, or "left-over" variation, that is unexplainable by any of the other components. Assuming independence between each component, the variance of the value of $P$ can be decomposed according to Eq. 2.10.

$$P = \mu_{LOT} + \mu_{WAFER} + \mu_{DIE} + \mu_{X,Y} + \epsilon \tag{2.9}$$

$$\sigma^2(P) = \sigma^2_{L2L}(P) + \sigma^2_{W2W}(P) + \sigma^2_{D2D}(P) + \sigma^2_{WID}(P) + \sigma^2_{\epsilon}(P) \tag{2.10}$$

Technically, the left side of Figure 2-9 can be extended to include fab-to-fab (in the case of multiple fabs manufacturing the same product) and tool-to-tool components as well. However, without loss of generality, these components can be lumped into the lot-to-lot component. For the purposes of this thesis, we will distinguish primarily between die-to-die, within-die and random variation for the following reasons:

1. The chips fabricated in this work all belong to the same wafer (or two wafers at most) so there is no statistical significance in wafer-to-wafer or lot-to-lot decomposition.

2. Circuit designers and architects are mainly concerned with ensuring that the power and performance of sellable products, typically single die, fall within defined constraints. In this context, other variation components can be included as additional mean shifts in the die-to-die component.

To include the effects of the other variation components, when we refer to die-to-die variation, we will, unless otherwise stated, also include the lot-to-lot and wafer-to-wafer components, if they exist and can be extracted.

As the significance of process variation increases, growing effort has been made to decompose spatial variation at smaller and smaller length scales. Indeed, at the

extreme, LER and RDF provide examples of within-transistor variation length scales. The within-die component of variation includes both systematic (repeatable within-die variation pattern from one die to another) and random components (e.g., $\mu_{X,Y}$ and $\epsilon$), with the random component being further decomposed into spatially correlated (where the distance between two devices or circuits determines how correlated those devices or circuits are, as shown in Figure 2-10(a) and Eq. 2.11) and uncorrelated components. This manner of distinction and classification is critical for the following reasons. Firstly, understanding whether systematic, repeatable variation patterns exist allows designers to tailor solutions that are specific to the existing variation pattern (perhaps easier and requiring less overhead), or must instead be broad enough in scope to allow for any possible pattern. Secondly, knowledge of spatial distance-dependent correlation directly impacts the cumulative amount of variation expected in large circuits. To mathematically illustrate this, consider the variance of the sum of identical, normally distributed random variables: when uncorrelated, the variance of the sum is $\sigma_{TOT}^2 \approx \frac{\sigma_{IND}^2}{n}$. However, when perfectly correlated ($\rho = 1$), the variance is $\sigma_{TOT}^2 \approx \sigma_{IND}^2$, for large $n$, which is a potentially large difference as $n$ increases.

In this context, we can formulate a two-axis system to decompose within-die variation as shown in Figure 2-10(b). As examples, the axes are pre-populated with two of the process steps discussed in Section 2.1, CMP and Lithography. Variation in CMP is largely a result of feature sizes, neighboring features, and pattern-dependencies [27], leading to both highly repeatable (systematic) within-die variation patterns from one die to another, as well as high spatial correlation between devices on a single die. Lithography variation is similar in that it is highly systematic due to mask errors, lens aberrations and other such effects that influence each die in the same manner, but different from CMP variation as the spatial influence of sub-wavelength lithographic systems with OPC is on the order of $1 - 2\mu m$ [44].

$$\rho_{x,x\pm d}(d) = \frac{Cov(x, x \pm d)}{\sigma_x^2} \tag{2.11}$$

Throughout this thesis, we will update the axes in Figure 2-10 with the results of

(a) Spatial (distance-dependent) correlation

(b) Spatial variation decomposition axes

Figure 2-10: Decomposition of within-die spatial variation

measured data from fabricated test-chips.

## 2.4 Impact to Modeling and Design of Circuits and Systems

Given the increasing challenge that variation poses to further CMOS scaling, there has been a greater drive to characterize, analyze and better understand the sources of variation as well as their circuit implications.

Complete understanding of the physical processes resulting in certain types of variation is theoretically possible. For example, CMP of interconnect metals is close to a point where it can be modeled accurately enough to substantially reduce layout-induced variations resulting from this unit process [45, 28]. Optical Proximity Correction (OPC) and Sub-Resolution Assists Features (SRAF) are used in the lithographic portions of the manufacturing process to reduce uncertainty in the patterned silicon features. These tools are a result of fairly well understood and modeled optical phenomena that occur when patterning features smaller than the wavelength of light used in the lithographic system. However, the decreasing physical dimensions necessary to satiate the industry's desire to further Moore's Law continue to outpace even these

"assists."

More often it is computationally or otherwise practically prohibitive to physically understand and model individual variation sources. In such cases, statistical modeling is employed to infer device or circuit behavior in the face of hypothesized or measured process variation. The simplest form of statistical modeling involves large numbers of Monte Carlo simulations to characterize performance to an acceptable confidence interval. More recently, faster statistical techniques have been formulated due to the increasing attention being paid to variability. However, many of these techniques remain primarily in the academic context because of their relative immaturity and inability to robustly handle certain phenomena such as parameter and device correlation.

To cope with such uncertainty and lack of appropriate modeling, designers often use feedback mechanisms to correct circuit behavior dynamically regardless of the underlying variation mechanisms. A simple example is the Phase-Locked Loop (PLL) used in countless analog and digital designs. This circuit utilizes a closed-loop feedback structure to create a stable clock of known frequency that is phase-locked (and often orders of magnitude faster in frequency) with an input clock, regardless of the underlying variations in the circuit generating the faster clock. Such mechanisms require an appropriate variation sensor and the ability to dynamically modify the behavior of the underlying circuit. The PLL uses a phase detector as the sensor, the output of which drives a charge pump which generates a voltage to control the Voltage-Controlled Oscillator (VCO).

However, in many feedback loops, replica circuits, or copies of the actual circuit, are the variation sensors. The behavior of these replica circuits is assumed to be the same as that of the actual circuit whose behavior will be modified by the feedback loop, even if the distance between the circuits is significant. Newer technology nodes and even the region of circuit operation (e.g., sub-threshold operation) make this assumption less valid. As a result, it is critical to understand and quantify the correlation between the behavior of the actual circuit and that of the replica circuit. More generally, understanding the spatial correlation of arbitrary variation sensors

and the circuits they monitor, even if their circuit structures lack any commonality, is desirable. This understanding then allows designers to appropriately design and place sensors such that they provide the most impact in representing and mitigating the effects of process variation with the least overhead.

As the computing and technology industries move to the "multi-core" regime, where many small to medium size cores work in parallel to achieve computational throughput, understanding spatial correlation in device and circuit behavior becomes critical. Developing circuit and architectural techniques that take this understanding into consideration to mitigate the effects of process variation are equally important in allowing computational ability to scale further. Thus, the primary objective of this thesis is to develop test structures to quantify spatial variation at both the device and circuit level, and to subsequently identify circuit and architectural mechanisms of compensating or coping with the impact of process variation.

## 2.5 Background & Related Work

The body of literature relating to variation in the manufacturing process is extensive and ranges from completely process related, to circuit and system design in the face of device and interconnect uncertainty. The pace of research in the field has accelerated in recent years due to its increasing significance [17]. However, this field is hardly new: as early as 1969, researchers were exploring correlations between non-uniformity in diffusion processes and the resultant electrical characteristics of bipolar transistors [46]. A few years later, studies on threshold voltage variation as a result of process variation began to appear [47, 48].

Nevertheless, most early work in the field relates to manufacturing process control to ensure acceptable process windows and yields. This work remained primarily the focus of and limited to manufacturing lines with designers rarely affected, owing to small relative variances in design parameters. In the early 1990's exploration of the impact of process variation on sensitive circuits such as high-speed analog circuits, SRAMs and clock distribution networks, and circuit techniques to mitigate this im-

43

pact began to appear [49]. To this day, these circuits remain some of the most critical and sensitive to variation of any type. However, variation remained largely out of the consciousness of most designers and researchers until approximately the $180nm$ CMOS technology node. At this node, the patterns lithographically "printed" on silicon were for the first time smaller than the wavelength of light ($193nm$) patterning them, resulting in sub-wavelength optical phenomena introducing undesirable effects. Clever process tricks, as well as good modeling of these effects, were able to stave off large parameter uncertainties for another couple of technology nodes, at which point variation was no longer simply wafer-to-wafer or die-to-die shifts. Instead, the variation problem began to significantly affect circuits within-die, meaning two identical circuits on the same die behaved differently [50]. This resulted in a flurry of research on the underlying causes of variation, efficient modeling methods to enable designers to study the effects during design, changes or additions to the manufacturing process to reduce absolute uncertainties, and circuit techniques to moderate the impact to sensitive circuits. As scaling continues, researchers now look for different transistor structures which not only improve electrical performance relative to typical bulk CMOS, but also improve the structure's inherent sensitivity to process variations [51]. Furthermore, statistical analysis and optimization techniques in the context of process variation are now at the forefront of research in this area [52].

The remainder of this section discusses specific work at the device, circuit and architecture levels focused on both characterization and mitigation of variation at those levels, while emphasizing spatial decomposition of the variation.

## 2.5.1 Characterization

Table 2.4 shows the most significant characterization methodologies in use at each of the levels of the hierarchy, with details in the sections below.

| Hierarchy Level | Methodology |
|---|---|
| Devices | Direct probing, pad multiplexing [53], large array device test-structures [54, 55, 56, 57, 58, 59] |
| Circuits | Ring oscillators [60, 61, 62, 63], SRAMs [58, 64] |
| Architecture & Systems | Simulation (FPGA [65], Microprocessor [66, 67, 68]) |

Table 2.4: Characterization methodologies for each level of the design hierarchy.

## Devices

A large body of work exists in studying transistor characteristics by measuring the current versus voltage curves of isolated transistors. These test setups involve probing individual transistors using four probe pads, one for each terminal of the transistor. More recent designs enable many individual transistor characteristics to be studied more efficiently without the use of dedicated per-device probe pads [53]. In mature technology nodes such studies not only provided I-V data for use in creating pseudo-physical transistor models for simulation, they also provided sufficient data to characterize the variation in that technology node. However, as dimensions shrink beyond the wavelength of light patterning them, non-intuitive physical phenomena result in new types of variations including interactions between the designed layout and unit process modules, and even interactions between devices laid out in close physical proximity. Accordingly, it is no longer sufficient to only study individual devices in isolation.

While the field is large and broad in scope, there has been little published relating to the spatial decomposition of process variations at the die level. This is not entirely surprising as companies often claim that such information is proprietary and critical to their success. The first published work relating variation to spacing information was that of the local mismatch model proposed by Pelgrom et al., which related device variation to the area of the device and space between the devices being studied [41].

Characterization of spatial variation requires large numbers of replicated test structures to capture variation statistics as a function of separation distance. Realizing this, Kibarian et al. attempted to create models for correlating test-structure data from a $1.2\mu m$ process with spatial information as early as 1990 [69], but no

(a) Schematic of addressable, large array test structure from [54].

(b) $V_T$ spatial correlation results from [55].

Figure 2-11: IBM's efforts in spatial variation characterization.

subsequent analysis with advanced technology nodes was published for several years. Within the past four years there has been substantial progress in developing arrayed test structures capable of such characterization. Perhaps the most significant work in this area concerns channel length variation [70, 16], which showed that after subtracting layout-dependent systematic variation, the remaining within-die variation is very weakly or not at all spatially correlated. Increasing attention is being paid to $V_T$, primarily in the context of SRAM cells which are particularly sensitive to $V_T$ variation. Agarwal et al., Rao et al., and Mukhopadhyay et al. at IBM [54, 55, 56, 57], Fischer et al. [58], and Wang et al. [59] all have developed efficient measurement and characterization of arrayed structures, but only a subset of these present data regarding spatial correlation analysis of the measured results. Agarwal et al. present spatial correlation analysis of SRAM-sized devices in a $65nm$ SOI process, concluding that little spatial correlation in $V_T$ within each die exists at $65nm$. Fischer et al. also provide $V_T$ variation data and autocorrelation analysis of 1-M SRAM cells at both the $90nm$ and $65nm$ nodes, with the results again showing no spatial correlation. Furthermore, both the aforementioned works make little to no effort to isolate threshold voltage variation from variation in channel length or other variation sources.

46

## Circuits

Test-structures that are ideally able to study individual device parameters as well as the behavior of the device within its circuit context have become necessary, as isolated transistors no longer truly represent "dense" transistors or those that are near each other in a circuit context. Making this difficult is the fact that many device parameters are coupled with other device parameters in a manner that does not allow them to be decoupled easily. Furthermore, multiple devices and device parameters often couple to a single circuit metric, impeding the decoupling of circuit behavior from that of multiple parameters within a device or even multiple devices.

Nevertheless, a number of groups, including our own, have studied small, simple circuits to gain understanding of both device and circuit behavior in the face of variation. Often, ring-oscillator frequency is used to characterize variation at both within-die and die-to-die levels [60, 61, 62, 63], due to the simplicity of the circuit and ease of frequency measurement. While ring-oscillator based techniques enable fast, easy test setups, this comes at a price: isolation of individual parameters for variability study is challenging due to amalgamation of the variation of many transistors into a single parameter (i.e., the frequency of ring operation). In some cases, it is possible to isolate the effects of critical device and process parameters. The authors of [63] are able to isolate a single, individual device parameter, $V_T$, in a small number of transistors within the ring, by including these transistors in pass-gate configurations between each inverter stage of the ring. By using short rings, they are also able to limit the averaging occurring due to parameter lumping. As this example illustrates, careful design is necessitated to decouple device variation from overall circuit variation.

Another small circuit commonly investigated is the six transistor SRAM cell due to its small size and aforementioned susceptibility to variation. Rather than look at the individual devices within these cells, Guo et al. develop a methodology to investigate variability in the read/write margin of the entire SRAM cell in large arrays [64]. Analysis of the measured data, other than observation of systematic processing effects and expected mismatch, is not provided.

47

## Architecture & Systems

Not surprisingly, most characterization of variation at the architecture or system level has occurred in the simulation domain, as characterizing and comparing architectures is impossible after a particular architecture has been chosen and implemented. As this is an emerging area of research, to date, most work has concentrated on individual metrics such as full-chip timing or leakage.

In the past, analog circuits and systems have been more susceptible to variation and have received more attention even at the system level. Conroy et al. use a statistical modeling framework to investigate various segmentation architectures for high-resolution digital-to-analog converters (DACs) [71]. Realizing that spatial correlation can drastically affect the linearity of DACs, they implement an interactive Monte-Carlo based statistical modeling framework capable of handling spatially correlated variations.

In the digital realm, FPGA architecture evaluations were among the first to consider process variations at this level [65]. Timing and leakage variation models are developed with the aid of Monte-Carlo simulations. However, all parameter distributions are assumed to be Gaussian, with no consideration of spatially correlated parameter variation. The assumption of Gaussian parameter distributions is mathematically and practically easier to work with and, in some cases, such as SRAMs, this is often an accurate assumption. But, in many circuits this type of analysis often leads to over-engineered solutions, as consideration of spatial correlation often leads to reduced bounds on circuit variability.

More recently, Intel has published work that details the impact of systematic variation on full-chip timing, indicating that such variation makes critical paths even worse than purely random variation would [66]. The authors of the study provide several near- and long-term solutions proposed in the literature, but make no mention if any of these techniques were used in their own products. Work has also been done characterizing power variability in multicore processor architectures, but only random variability is taken into account [68]. The authors of [67] show that correlated

variation, in particular correlated channel length variation, impacts the frequency of cores within a multicore processor.

## 2.5.2 Mitigation

Mitigation of process variation is best divided into process-oriented versus design-oriented approaches. Targeting the process itself involves altering process modules and/or flows or device design, directly impacting variation at or close to the source. Design-oriented strategies — through design rules, layout modifications, or robust circuit and system design — offer an indirect route to variation mitigation. The following table summarizes the more significant mitigation strategies and related work in those areas for each level of the hierarchy. Once again, details of each are provided in the following sections.

| Hierarchy Level | Methodology |
| --- | --- |
| Devices | Statistical process control [72, 73, 74, 29], Resolution enhancement [75], materials improvement [10] |
| Circuits | Feedback, Circuit redesign for robustness [76, 77], Adaptive scaling [78, 79, 80, 81] |
| Architecture & Systems | Redundancy [82, 83], Voltage/Frequency scaling and islands [84, 78, 80, 67] |

Table 2.5: Mitigation strategies for each level of the design hierarchy.

**Devices**

At the device level, it is more relevant to discuss process-oriented variation reduction approaches. Statistical process control made its debut roughly 20 years ago as a means to reduce "back-end" testing and assembly costs in IC fabrication lines by increasing yields in the "front-end" steps described above [85]. In-situ measurements coupled with inferential signals from individual process modules are utilized to dynamically change process conditions and recipes to ultimately benefit yield. In the ensuing years, statistical process control has pervaded nearly every process module, including oxidation [72], photoresist application [73], etch [74], and CMP processes [29].

49

In addition to feed-back/forward process control, continuous improvements have been made to process steps and in the use of materials to improve both device performance and variability. The resolution enhancement techniques described in Section 2.1.1, immersion lithography [75], other patterning improvements and introduction of high-k/metal-gates to high-volume CMOS manufacturing [10] (as shown in Figure 2-12) all provide examples of techniques meant to improve manufacturability of high-performance transistors. Whole process steps, such as CMP, have been added for this sole purpose.



Figure 2-12: Improvement in $\sigma_{V_T}$ in Intel's process due to oxide scaling and new materials, though there is no mention of how $\frac{\sigma_{V_T}}{\mu_{V_T}}$ scales. $C_2$ is analogous to $A_{V_T}$ found in Eq. 2.7. Reproduced from [11].

Common to modern semiconductor manufacturing are combination "process-design" mitigation techniques, where process and design teams actively engage each other to trade-off ease-of-design, manufacturability and performance. This process is best implemented where design and manufacturing occur in the same organization as evidenced by Intel [10], but Design For Manufacturability (DFM) kits are increasingly made available by foundries as well.

## Circuits

Digital circuits and systems, due to both their bi-stable nature and high gain in the switching region which results in relatively large noise margins, are inherently more robust to variation. Analog circuits, however, are considerably more sensitive

to process variation and have employed both layout and design-oriented mitigation strategies — including common-centroid layout, use of larger than minimum-size devices and feedback — for many years to cope with process variation. Among the first circuits to employ design-oriented mitigation was the PLL, which made use of feedback to decrease sensitivity to component values [86].

PLLs are not limited to analog systems though; they are critical components of high-performance digital systems, providing the clock(s) necessary for timing of digital circuits. Being among the most sensitive and critical signals in digital systems, clocks and clocking networks were subject to variation mitigation strategies long before other portions of digital systems. Attempts to reduce clock skew and jitter have produced a myriad number of mitigation schemes such as H-Tree, X-Tree and grid clock distribution, active deskewing [87, 83], and even some proposed exotic solutions such as use of optical clock networks [88, 89]. Also a subject of much research in variation mitigation are SRAM circuits in cache memories, which are particularly susceptible to variation due to their small physical size, leading to SRAM cell [90] and sense amplifier [76, 77] redesign as well as other so-called "peripheral assists" to mitigate the impacts of variation.

With relative process tolerances continuing to degrade, digital paths have become the focus of multiple variation mitigation techniques. Some of these include replica critical path monitor circuits around which adaptive feedback schemes can be built, such as Adaptive Body Biasing (ABB) [78], Adaptive Voltage Scaling [79, 80, 81], and detection and correction of timing errors within microprocessor pipeline stages [91, 12] (Figure 2-13). Portable electronics and sensors have motivated the need for low energy operation in the sub-threshold operating regime, in which transistor leakage currents are exponentially affected by variations in $V_T$ and which poses unique, but surmountable challenges to circuit design. With appropriate gate library modeling and design and variation-aware timing methodologies, sub-threshold circuits can be designed to be as robust as their above-threshold counterparts [92].

Understanding of spatial components of variation have also guided some mitigation work. In particular, Friedberg et al. claim such understanding can provide up to

Figure 2-13: RazorII latch capable of detecting timing errors and restarting the pipeline [12].

a 4x improvement in digital circuit delay variability if appropriate process control mechanisms are put in place to take advantage of spatial correlation data [70]. Lastly, statistical static timing, in which timing of digital paths is described by statistical distributions, has garnered attention in aiding designers assess the susceptibility of critical paths to variation and make appropriate changes to ensure high yield.

## Architectures & Systems

Architectures and systems are often too large and complicated to apply global techniques to mitigate variation. Often the most susceptible building blocks are the ones to which mitigation techniques are applied. For example, the susceptibility of SRAM caches to variation, noted in Section 2.5.2, has led to the development of "variation-aware" caches which dynamically resize as faulty bitcells are detected [82] or where whole cache lines are replaced with redundant lines [83]. Variation-aware synthesis techniques have also enabled mitigation in on-chip bus architectures for systems-on-a-chip, resulting in decreased power consumption as well [93].

Process variation has, however, reached the point at which system architects can no longer avoid considering its impact. Marculescu et al. include process variability in the context of heterogeneous blocks in an embedded application that together must meet some latency constraint [84]. Each block is a voltage/frequency island and the optimal voltages and frequencies for islands are solved for. Humenay et

52

al. evaluate Adaptive-Voltage Supply (AVS) and Adaptive Body-Biasing (ABB) as potential circuit-level solutions in the context of core-to-core frequency variations in multi-core processors, but conclude that both of these techniques involve tradeoffs between static and dynamic power, making them less desirable potential solutions [67]. Similarly, an adaptive FPGA architecture utilizing ABB is proposed in [94] leading to 3.45X reduction in timing variability and 3X reduction in leakage power. Donald et al. also explore core-to-core variation in frequency, and propose allowing the system to turn off cores if the additional power consumed by the core is higher than a proposed metric [95]. Most recently, Liang et al. use voltage interpolation at the gate or pipeline level of microprocessors to mitigate process variability [96] and potentially reduce energy consumption considerably.

### 2.5.3 Modeling & Simulation

One community that has put forth a large body of work involving spatial variations is the CAD/modeling community. Statistical modeling of power as well as sophisticated statistical timing analysis including the effects of spatial correlation have been of great interest in recent years [70, 97, 98, 99, 100, 101, 102, 103, 104, 105]. In particular, many of these attempt to formulate mathematically efficient methods, in the context of spatially correlated variation, to predict the probability density function (PDF) and associated cumulative density function (CDF) of various metrics (most often timing or leakage power) for a given design.

Either due to the relative immaturity of these analysis and modeling techniques or lack of sufficient data, there has been no published work utilizing these tools to gain insight into techniques for mitigation of spatially correlated process variation.

# Chapter 3

# Device Parameter Variation: Threshold Voltage

The previous chapter described the many sources of variation in semiconductors and highlighted the need to characterize the spatial aspect of this variation. We begin by studying variation at the device level, as devices are the fundamental building blocks of circuits and systems and most modeling and simulation begins at this level. This chapter first motivates the need to study variation in two important transistor parameters, channel length ($\triangle L$) and threshold voltage ($V_T$) in Section 3.1, followed by theory and implementation details of a test-chip fabricated to characterize $V_T$ variation in Section 3.2. Comprehensive analysis, focused on characterizing the spatial components of the observed variation data and which shows little to no spatial component of $V_T$ variation, is presented in Section 3.3. Lastly, Section 3.4 explores the resulting implications for both the modeling and simulation as well as the circuit design communities.

## 3.1 Motivation

Section 2.2 showed that the performance of digital circuits is determined, to first order, by five parameters: mobility ($\mu$), oxide capacitance ($C_{ox}$), transistor dimensions ($W$, $L$), and threshold voltage ($V_T$). Modeling variation in these parameters requires

relating each to the more fundamental parameters $N_{ch}$, $t_{ox}$, $W$, and $L$. However, practically measuring variability in many fundamental parameters (e.g., $N_{ch}$) is often extremely difficult. Furthermore, both oxide thickness (capacitance) and mobility are relatively well controlled, with variation typically much smaller than the 10% channel length variation and 30% threshold voltage variation. As such, variation in channel length and threshold voltage is more critical and relevant to circuit designers. Many modeling and simulation programs, such as Hspice, facilitate modeling variability of these critical parameters through the use of parameters like *delvto* and direct manipulation of $L$ in the netlist [106].

Friedberg et al. have undertaken extensive work in characterizing and decomposing variation, *including spatial decomposition and characterization*, in lithographically critical dimensions, particularly channel length [16]. Careful analysis revealed that within-die "spatial correlation is virtually entirely an artifact of systematic variation." Note here the identical distinction being made between separation-distance dependent correlation and systematic versus random variation as that shown in Figure 2-10. These results provide further justification for our placement of lithography as systematic but uncorrelated on that figure.

The aforementioned criticality of $V_T$ to circuit designers requires analogous analysis and decomposition of threshold voltage variation for more complete understanding of process variation at the circuit level. Since $V_T$ is not a fundamental parameter but is instead dependent on several fundamental parameters, careful attention must be paid in its extraction to be of value in modeling and simulation. In particular, the following section describes how measurement of sub-threshold currents of transistors allows for isolation and extraction of the intrinsic threshold voltage, $V_{T_o}$. When coupled with variation in channel length, a complete picture of both channel length and threshold voltage variation begins to emerge.

## 3.2 Theory & Enabling Circuits

In this section, we describe the device operation fundamentals and analytic basis used to extract $V_T$ variation from leakage current measurements, as well as the overall test chip architecture and circuits.

### 3.2.1 Extraction of $\triangle V_{T_o}$

The intrinsic or ideal threshold voltage, $V_{T_o}$, of a MOSFET is defined by Eq. 3.1, where $V_{FB}$ is the flat-band voltage, $\phi_{F_p}$ is the Fermi potential of the substrate, $\gamma$ is the body-factor and $V_{SB}$ is the source to body bias. Thus, $V_{T_o}$ is fundamentally dependent only on substrate doping, $N_{sub}$, and the oxide thickness, $t_{ox}$, through $\gamma$.

$$V_{T_o} = V_{FB} + 2\phi_{F_p} + \gamma\sqrt{2\phi_{F_p} + V_{SB}} \qquad (3.1)$$

However, due to short-channel effects, notably Drain-Induced Barrier Lowering (DIBL), in the deep sub-micron regime, the actual threshold voltage, $V_T$, becomes a function of not only $N_{sub}$ and $t_{ox}$, but also channel length, $L$, and source/drain junction depth, $x_j$ [5]. To quantify this dependence, $V_T$ is now defined as a summation of the ideal threshold voltage $(V_{T_o})$ and a shift $(\triangle V_T)$ due to short channel effects, as in [107].

Since $\triangle V_T$ formulated in [107] is primarily impacted by channel length (exponential dependence), we instead seek to characterize the variation and, in particular, spatial variation of $V_{T_o}$ to enable complete modeling and simulation of $V_T$ variation. This requires isolation of $V_{T_o}$ from other common sources of variation such as $\triangle L$, and can be achieved to a large degree in the sub-threshold regime of transistor operation.[1]

In this regime, the current through the transistor is given by Eq. 3.2, where $I_o$ is the drain current at $V_{GS} = V_{T_o}$, $\gamma'$ is the body-effect coefficient, $\eta$ is the DIBL coefficient, $n$ is the sub-threshold slope ideality parameter defined by Eq. 3.3, and

---

[1]From this point forward we will use $V_T$ and $V_{T_o}$ interchangeably indicative of characterization of the *ideal* MOSFET threshold voltage.

$V_{TH}$ is the thermal voltage.

$$I_D = I_o \cdot e^{\frac{V_{GS} - V_{T_o} - (\gamma' \cdot V_{SB}) + \eta \cdot V_{DS}}{n V_{TH}}} \cdot \left(1 - e^{\frac{-V_{DS}}{V_{TH}}}\right) \tag{3.2}$$

$$n = \frac{\triangle V_{GS}}{V_{TH} \cdot \triangle log(I_D) \cdot ln(10)} \tag{3.3}$$

The $\left(1 - e^{\frac{-V_{DS}}{V_{TH}}}\right)$ term in Eq. 3.2 is easily eliminated with $V_{DS} > 3 \cdot V_{TH}$, and $\gamma' \cdot V_{SB}$ can be eliminated by shorting the body and source of each device. Minimization of the DIBL component is achieved by reducing $V_{DS}$ to a few hundred millivolts — large enough to eliminate the $(1 - e)$ term, but still small enough to minimize the effect of DIBL. Assuming that $V_{GS}$, $V_{DS}$, $\eta$, and $n$ are identical for two arbitrary devices, taking the natural logarithm of the ratio of the currents of those two devices will result in a simple analytic equation for the $\triangle V_T$ of those devices[2], as shown in the following derivation:

$$ln\left(\frac{I_{D_2}}{I_{D_1}}\right) = ln\left(\frac{I_o \cdot e^{\frac{V_{GS} - V_{T_{o2}} + \eta \cdot V_{DS}}{n V_{TH}}}}{I_o \cdot e^{\frac{V_{GS} - V_{T_{o1}} + \eta \cdot V_{DS}}{n V_{TH}}}}\right) \tag{3.4}$$

$$ln\left(\frac{I_{D_2}}{I_{D_1}}\right) = ln\left(e^{\frac{V_{T_{o1}} - V_{T_{o2}}}{n V_{TH}}}\right) \tag{3.5}$$

$$\triangle V_{T_{1,2}} = n V_{TH} \cdot ln\left(\frac{I_{D_2}}{I_{D_1}}\right) \tag{3.6}$$

It is known that the subthreshold slope between devices can vary, and thus the assumption that $n$ is identical for the devices being compared is not valid in general. Taking this into account, the above equations can be reworked and shown to provide the relationship in Eq. 3.7 between the drain currents, threshold voltages and $n$.

$$n_1 V_{th} \cdot ln\left(\frac{I_{D_1}}{I_{D_2}}\right) - \frac{n_2 - n_1}{n_2} \cdot V_{GS} = \frac{n_1}{n_2} \cdot V_{T_2} - V_{T_1} \tag{3.7}$$

Since Eq. 3.7 provides no simple, closed-form solution for $\triangle V_T$, the values of $n_1$,

---

[2]This $\triangle V_T$ denotes the difference in ideal threshold voltage between two devices and is distinct from $\triangle V_T$ in [107] which describes the shift in a single device's threshold voltage due to short channel effects.

$n_2$, $V_{GS}$ and at least one device $V_T$ must be known to compute the other and thus a delta between the two. Eq. 3.3 shows that $n$ can easily be computed using two measurements of $I_D$ at differing values of $V_{GS}$. We also note that by using a small value for $V_{GS}$ we can minimize the contribution of the second term in Eq. 3.7 in two ways: a smaller $V_{GS}$ results in 1) a smaller multiplicand, and 2) as can be seen in Figure 3-1, the instantaneous value of $n$ tends to converge at extremely low $V_{GS}$ despite variation in $V_T$, allowing use of Eq. 3.6 rather than Eq. 3.7.

Ascertaining the value of one of the device threshold voltages is more difficult, but possible by finding the value of $V_{GS}$ where the sub-threshold current deviates from the ideal log-linear form. This known $V_T$ can then be used to compute the $V_T$ for every other DUT. By using Eq. 3.7, along with two additional measurements for each DUT (to compute $n$), a complete $V_T$ map of all test devices across the chip can be ascertained.



Figure 3-1: Convergence of $n$ at low $V_{GS}$ despite $V_T$ variation.

A more practical method is to compute an average value, $n_{avg} = 0.5 \cdot (n_1 + n_2)$, where each $n$ can be computed from sub-threshold leakage current measurements, and then use $n_{avg}$ in Eq. 3.6. This has the added benefit that post-processing of the current measurements to extract $\Delta V_T$ remains computationally efficient. Furthermore, the error associated with using this average value is small if $n_1 \approx n_2$.

## 3.2.2 Test-Structure Architecture & Circuits



Figure 3-2: Simplified $V_T$ variation architecture and circuits

Figure 3-2 is a simplified schematic showing the architecture and circuit blocks to measure $V_T$ variation. A dual-slope, integrating Analog-to-Digital Converter (ADC) is used to measure sub-threshold leakage currents, due to its suitability to accurately measure small currents despite long conversion times. The resolution of the ADC in this design is configurable up to 13 bits, allowing a trade off between accuracy and measurement time. By externally setting $V_{DS_{ref}}$, the operational amplifier enforces a virtual ground at the input nodes, ensuring that each device connected is biased at the same $V_{DS}$. Amplifier gain and mismatch errors will introduce error into the value of $V_{DS}$ seen by the DUTs, but each DUT is affected in the same manner. Large input devices and a high-gain ($> 60dB$) ensure that this error is small nevertheless.

To measure currents of many devices efficiently, we use a single ADC that is multiplexed among all devices. Apart from the area efficiency achieved by using a single ADC, this ensures that any non-idealities in the ADC are common to all DUTs and therefore do not affect the results. We have chosen to use a hierarchical access scheme analagous to a memory, with rows, columns and sections. Each bank contains 128 PMOS and 128 NMOS DUTs organized in columns, as shown in Figure 3-3, and each section contains 90 rows of banks. Bank enable pass gates steer only the current of the selected device to the measurement circuitry. The test-chip contains 540 banks organized into 6 sections for a total of 540 banks x 128 columns = 69,120 DUTs of

each type in a 2mm x 2mm array. All device lengths are minimum length for this technology ($0.18\mu m$). The lower portion of this array contains banks with random designed device width ranging from $0.28\mu m$, the minimum allowable, to $3\mu m$, while in the upper half of the array, each row contains the same device width.



Figure 3-3: Simplified schematic of individual bank

While the row and column access transistors introduce resistance and variation, Hspice simulations show that a $\pm10\%$ variation in either $L$ or $V_T$ of the access transistors has $< 0.5\%$ effect (Section 3.2.3) on $I_{DS}$ of the DUT being accessed. Despite column access transistors being turned completely off for all other DUTs, a finite leakage current, $I_{leak}$, through the row and column access transistors and the "off" DUTs adds to the $I_{DS}$ of the DUT being accessed. When $I_{DS}$ of the accessed DUT is large, corresponding to a larger $V_{GS}$, $I_{leak}$ is a negligible component and can safely be ignored. However, as described in Section 3.2.1, it is desirable to set $V_{GS}$ as low as possible to benefit from the convergence of $n$ at low $V_{GS}$. At gate biases below 0.25V, $I_{DS}$ reduces to nanoamps or smaller, so that even small drain-source leakage currents and drain/source-body junction currents accumulate over the "off" DUTs and their access transistors.

Each of the bank enable pass gates are high-$V_T$ devices to minimize these parasitic leakage currents emanating from unaccessed devices. However, even with very small parasitic leakage currents from each of these pass gates, the large number of gates result in these parasitic currents summing to a current large enough to interfere with the current being measured, degrading the signal-to-noise ratio. An active current

61

subtraction scheme was devised and implemented on-chip as shown in Figure 3-2. Two (source and sink) 8-bit thermometer-code DACs [108] with digital control logic are used to actively add or subtract current equivalent to $-I_{leak}$. The digital logic implements a binary search algorithm that uses the output of the ADC to converge upon the correct DAC input value, acting as an auto-zeroing mechanism. For example, when trying to measure the first NMOS DUT in a bank, the auto-zeroing is first run when all DUTs in the bank are off. If the ADC output is anything but 0 after the first auto-zeroing measurement, the digital logic will completely turn on one of the two DACs shown in Figure 3-2 in response to the direction of $I_{leak}$. If $I_{leak}$ is being drawn from the measurement circuitry to ground, the algorithm turns on the source DAC (top of Figure 3-2) in order to "source" $I_{leak}$ and remove its effect from the measurement. Analogously, the algorithm will turn on the sink DAC (bottom of Figure 3-2) to "sink" an $I_{leak}$ flowing from $V_{DD}$ to the measurement circuitry. Subsequent measurements are used to refine the DAC control word in a logarithmic fashion. Auto-zeroing is performed once for each bank being tested at a specific gate bias. Due to the discrete nature of a DAC as well as limited resolution, the auto-zeroing will not be perfect, and residual $I_{leak}$ is treated as an offset and subtracted from DUT current measurements.

The test-chip was implemented on a National Semiconductor $0.18\mu m$ bulk CMOS process. Figure 3-4 is a die photo showing the 3.2mm x 2.7mm test-chip, of which 2mm x 2mm is the dense DUT array.



Figure 3-4: Test-chip die photo. The DUT array is shown at left, with the ADC and digital control and calibtration blocks at right.

### 3.2.3 Simulation Results of $\triangle V_T$ Isolation

We next present simulation results showing that the above theory and circuits are effective in isolating and extracting $\triangle V_T$ even in the presence of other types of variation, particularly channel length variation. Furthermore, simulations show that the multiplexing circuitry contributes a negligible error in the measured current.

**$V_T$ Isolation**

Simulations were performed in which a single DUT, with row and column access transistors, charges an integrating capacitor. Since the operational amplifier forces a virtual ground at the inputs, $V_{DS}$ remains constant as the capacitor is being charged, and $I_{DS}$ can then be found using $I_{DS} = C \cdot \triangle V / \triangle t$. Simulations were performed with DUT $V_T$ and channel length being varied by $+/-10\%$ and $V_{DS}$ being varied across the allowable range, determined by the output stage of the operational amplifier. The range in this design is 0.3V - 1.5V.

The plots in Figure 3-5(a) show the simulation results. The top plot in the figure varies DUT $V_T$ and $V_{DS}$, while the bottom plot varies DUT channel length and $V_{DS}$. Both plots show the relative change in current from the nominal $V_T$ or $L$ at a given value of $V_{DS}$. The plots clearly show that the arrangement detailed above is more sensitive to changes in $I_{DS}$ as a result of $V_T$ variation rather than $L$ variation, especially at low $V_{DS}$. These results are consistent with the theory outlined in Section 3.2.1.

To quantify these results further, the sensitivity of $I_{DS}$ to either $\triangle V_T$ or $\triangle L$ can be computed by taking the derivative with respect to $\triangle V_T$ and $\triangle L$, respectively. Taking the ratio of these derivatives gives the relative sensitivity of the circuit to $\triangle V_T$ and $\triangle L$. Since it is clear from Figure 3-5(a) that the circuit is least sensitive to $\triangle L$ at low values of $V_{DS}$, these derivatives are only calculated for the lowest $V_{DS}$ value allowable, 0.3V. Figure 3-5(b) plots the ratio of $\frac{\triangle I_D}{\triangle L}$ to $\frac{\triangle I_D}{\triangle V_T}$ for $V_{GS}$ ranging from 0.35V to 0.5V. Lower values of $V_{GS}$ are not plotted, as a trend in decreasing sensitivity to $\triangle L$ with larger values of $V_{GS}$ is evident from the figure. However, Section 3.2.1 discussed employing lower values of $V_{GS}$, where the value of $n$ converges

(a) Percent change in current from nominal value at a given $V_{DS}$ when varying $V_T$ (top plot) and channel length, $L$ (bottom plot)

(b) Circuit Sensitivity Ratio

Figure 3-5: Current sensitivity to variation in $V_T$, $L$

despite variation. The results of sensitivity analysis imply that simply measuring the value of $n$, as discussed in Section 3.2.1, and using a $V_{GS}$ near the nominal $V_T$ for the process provide more benefit in extracting $\triangle V_T$ than attempting measurements at extremely low gate biases. Furthermore, low gate bias values increase the resolution and dynamic range requirements of the ADC, and even more attention must be paid to preventing $I_{leak}$ currents from unselected DUTs.

For all values of $V_{GS}$ plotted, the sensitivity ratio through the majority of the variation range is below 0.1, meaning the circuit is at least 10X more sensitive to $V_T$ variation than to $L$ variation. This is particularly true for $|V_{GS}| = 0.35V$ where all but the endpoints remain $< 0.1$. Furthermore, sensitivity to $V_T$ variation peaks in the vicinity of the nominal values of $V_T$ and $L$. Since variation in these parameters are typically normally distributed about the nominal value, the majority of variation measured will be in the high-$V_T$-selectivity region of operation, giving high confidence that the measured $I_{DS}$ variation is primarily a result of threshold voltage variation.

Simulations were performed where a transistor was subjected to variation ($V_T$ or $L$ or both) and the ADC outputs used to determine the amount of variation with both Eq. 3.6 and Eq. 3.7. All simulations were done with $|V_{GS}| = 0.35V$ and $|V_{DS}| = 0.3V$ and the ADC resolution set to 10 bits. Simulations were also done at $|V_{GS}| = 0.345V$, the results of which are used in conjunction with Eq. 3.3 to calculate $n$. Table 3.2.3

contains the results of these simulations showing the ability of the circuit to measure the known variations. It should be noted that the simulations where only $V_T$ is varied result in approximately a 10% error in the extracted deltas, primarily a result of the inherent inaccuracy in using Eq. 3.6 which does not account for $n$ varying simultaneously. This alone would indicate that this test-structure can resolve deltas of approximately 1% of the nominal $V_T$. However, since the sensitivity of the circuit to $V_T$ variation is not infinite, resolution is reduced to approximately 2% of nominal $V_T$.

| Variation Type | Extracted $\triangle V_T$ |
|:---:|:---:|
| $+10\%V_T$ | $+10.9\%$ |
| $-10\%V_T$ | $-11.0\%$ |
| $+10\%V_T, +10\%L$ | $+11.9\%$ |
| $+10\%V_T, -10\%L$ | $+9.3\%$ |
| $-10\%V_T, -10\%L$ | $-12.0\%$ |
| $-10\%V_T, +10\%L$ | $-11.3\%$ |
| $+3\%V_T, +3\%L$ | $+3.6\%$ |
| $-3\%V_T, -3\%L$ | $-3.1\%$ |

Table 3.1: Extracted $V_T$ variation vs. subjected variation

## Access Transistor and Resistance Effects

In order to hierarchically access a large number of DUTs within an array, each DUT requires row and column access transistors, and each bank of DUTs requires a pass-gate. These devices introduce additional resistance, potentially lowering the $V_{DS}$ and corresponding $I_D$ of the DUT due to the finite $R_{out}$ of the devices. However, operation in the sub-threshold regime produces small currents which are not perturbed substantially by even fairly large resistances. Simulations were carried out to evaluate this impact. The test circuit was used in these simulations with and without row, column and bank access transistors. Table 3.2.3 shows that the impact of these transistors and variation within them is less than 0.5% of the simulated current without any access transistors.

Another possible source of inaccuracy in implementing a large array is variation

| Test Scheme | Relative Difference |
|---|---|
| DUT w/o access transistors | 0.00% |
| DUT w/access transistors | $-0.43\%$ |
| DUT with $-10\%\triangle L$ in access transistors | 0.26% |
| DUT with $-10\%\triangle V_T$ in access transistors | 0.23% |

Table 3.2: Simulated current differences due to inclusion and variation within access transistors ($|V_{GS}| = 0.35V$ and $|V_{DS}| = 0.3V$).

in the distance current must travel to the measurement circuitry, resulting in different resistances seen by each DUT to the ADC. However, simulations show that even with a $1k\Omega$ difference in resistance, the relative current difference is only 0.1%. Furthermore, the path from each DUT bank to the ADC is implemented as a dense metal grid spanning multiple metal layers to provide the lowest possible resistance. Process data and simulations indicate that a minimum width wire spanning 2mm has a resistance of approximately $500\Omega$. However, the grid is implemented with many 3X minimum-width wires spanning four metal layers, decreasing the overall resistance substantially. Since the resistance difference between any two DUTs cannot be more than the resistance of a single minimum-width wire spanning the entire array, we conclude that resistance variations in the grid will have negligible effect on measured currents.

## 3.3 Data Analysis

Next, we examine the measured currents from more than $50K$ devices on each of 36 chips. We first consider the ADC performance in measuring currents, followed by extraction of $\triangle V_T$ from our data. We then analyze device size and separation distance dependencies, and other within- and between-die spatial dependencies.

### 3.3.1 ADC Performance

Static performance of the ADC is first characterized, as this component is critical to our current measurements. Since the input to the ADC is a DC current, dynamic ADC

Figure 3-6: INL plot for a single ADC

performance is intentionally omitted as it has no effect on the measured currents. To alleviate noise concerns, measurements are repeated ten times and the average over the ten runs is used, as we assume white noise. In light of the limited input to the ADC, the primary metric we characterize is the Integral Non-Linearity (INL), as this gives the effective resolution of the ADC. Figure 3-6 is a plot of the INL versus ADC code for a single chip, when the ADC is configured for 10-bit resolution and $600nA$ full-scale current. At the high-end of the ADC range, the INL begins to degrade due to limited bandwidth of the operational amplifier. This limits the maximum full-scale current to $\approx 600nA$. However, a redesign of the operational amplifier could remove this limitation in future designs. The INL for the ADC shown in Figure 3-6 leads to an effective linearity of $10 - log_2 \left( MAX(INL) - MIN(INL) \right) = 7.81$ bits. Out of the 36 chips measured, the worst case effective linearity is $\approx 6$ bits, resulting in an effective resolution of $\frac{600nA}{2^6} = 9.375nA$. The minimum resolvable $\triangle V_T$ is then computed according to Eq. 3.6, where $I_2$ is the maximum current measurable by the ADC, $I_1$ is the current one ADC step below the maximum current and $n$ is conservatively estimated to be 1.5 for this process. This gives $\triangle V_{T_{MIN}} = 0.6mV$, or $0.14\%$ of the nominal $V_T$ for this process — below the 2% limit detailed in Section 3.2.3 — giving high confidence that ADC non-linearity contributes negligible error to the extraction of $\triangle V_T$.

67

Implementation limitations resulted in disabling of the auto-zeroing DACs intended to cancel the off-DUT leakage currents. However, each bank is still calibrated by first doing a current measurement with that bank's pass gate enabled but with no device enabled. Subtracting this measurement from the measured device current gives the true sub-threshold current of the enabled device, but limits the dynamic range of the ADC.

## 3.3.2 Current Measurements and Extracted $\triangle V_T$

Using the analysis in Section 3.2.1 and measured currents from DUTs within the array, we compute $\triangle V_T$ for each device with respect to a reference device in the corner of the array. The values of $n$ needed for the computation are extracted by measuring currents at different $V_{GS}$ biases ($0.275V \leq |V_{GS}| \leq 0.3V$ in $5mV$ increments) and computing a best-line fit on a semilog scale. The slope of this line relative to the ideal 60mV/decade gives $n$. A spatial plot of $n$ for one die is shown in Figure 3-7(a). From this plot, it is evident that $n$ does vary, although the magnitude is small, so using the average value of $n$ between two devices in computing the $\triangle V_T$ introduces only small errors.[3] Furthermore, due to the limits placed on the full-scale current by the degraded linearity at high current levels (Section 3.3.1), extracting an absolute $V_T$ to use in Eq. 3.7 is impossible, limiting calculation of $\triangle V_T$ to the formulation given by Eq. 3.6. This limitation can also be overcome in future designs.

Once $n$ has been computed from the measured currents for each device, we compute $\triangle V_T$. Figure 3-7(b) illustrates the $\triangle V_T$ from each device to the reference device in the bottom right corner of the array for an example chip. While it appears that there may be some correlation in the upper half of the array, this is only due to the systematic designed device width pattern from row to row in this section of the array, resulting in a shift in standard deviation but not correlation as will be shown in the following spatial analyses.

---

[3]The error is $< 10\%$, as this is the error when no attempt is made to account for $n$ varying as discussed in Section 3.2.3.

(a) Subthreshold ideality parameter $(n)$            (b) $\triangle V_T$

Figure 3-7: Spatial distributions of $n$ and $V_T$ for a single die

## 3.3.3 Pelgrom Modeling

To understand the effect of device size on standard deviation, we first show a Pelgrom plot of standard deviation of $V_T$ versus $\frac{1}{\sqrt{Area}}$ [41]. We note that $\sigma(V_T) = \sigma(\triangle V_T)$, where $\triangle V_T$ is with respect to our specified reference device. Based on Pelgrom's theory, we expect to see a linear relationship between $\sigma(V_T)$ and $\frac{1}{\sqrt{Area}}$. Figure 3-8(a) shows this linear relationship for a single chip, where $\sigma$ has been normalized relative to that for a device having W/L of $\frac{0.92\mu m}{0.18\mu m}$. Ideally, the best-fit line should pass through 0, indicating that devices of infinite area should have a standard deviation of 0. Deviation from this ideal in our data may be a result of ADC measurement resolution. Specifically, the data for many of the smallest devices in the array $(W < 0.8\mu m)$ is excluded as their extremely small currents are smaller than the resolution of the implemented ADC. In such cases, data is either non-existent (measured values being 0 as the currents were below the resolution of the ADC), or when data was present, gives unrealistic values of $n \gg 2$ due to the large ratio of ADC step size versus absolute current being measured at the bottom of the ADC range. Data is excluded by examining the computed values of $n$ and excluding any devices with a computed $n < 1.0$ or $n > 1.8$. The majority of the excluded data is for small device sizes in the tails of the distribution, artificially lowering the standard deviation for these device sizes. As a result, the data points on the far right of Figure 3-8(a) tend to be slightly

69

(a) $\sigma(V_T)$ vs $\frac{1}{\sqrt{Area}}$      (b) $\sigma(V_T)$ vs. Distance

Figure 3-8: Pelgrom fits of $\sigma(V_T)$

below the linear fit, and may also cause the fit line to no longer pass through 0. In all further plots and data analysis, the data for devices with $W < 0.8\mu m$ is excluded to ensure data integrity. We believe that all trends in the data remain applicable to the small device sizes in this technology.

The Pelgrom model includes two terms: 1) an area dependent term with proportionality $A_{V_{T0}}^2$, and 2) a distance dependent term with proportionality $S_{V_{T0}}^2$, as shown in Eq. 3.8. Figure 3-8(a) shows a clear dependence on device area, but Figure 3-8(b) shows no significant distance dependency. In Figure 3-8(b), all pairs of devices within a specific chip having the indicated separation distance $D$ are considered, and the standard deviation, normalized to a $\frac{0.92\mu m}{0.18\mu m}$ device, across all those pairs is plotted. Results show that placing devices nearer to each other does not decrease the variance between them. This will be discussed further in Section 3.4.

$$\sigma^2(V_{T0}) = \frac{A_{V_{T0}}^2}{WL} + S_{V_{T0}}^2 D^2 \qquad (3.8)$$

### 3.3.4 Intra-Die Spatial Correlation

Since $\sigma(V_T)$ depends on device size, separate spatial correlation analysis by individual device size is necessary. Intra-die spatial correlation analysis is performed by computing the correlation coefficient for all devices separated by some distance, $D$, where

(a) Intra-Die Spatial Correlation

(b) Die-to-Die Variation Pattern Correlation

Figure 3-9: Intra-die correlation vs. distance and die-to-die pattern correlation

$0 \leq D \leq 0.5 \cdot \sqrt{Array_X^2 + Array_Y^2}$ as shown in Eq. 3.9, where $DVT$ is the $\triangle V_T$ between a given device and the reference device. These correlation coefficients are then plotted versus separation distance in Figure 3-9(a) for a device size of $\frac{0.92 \mu m}{0.18 \mu m}$. The $\pm 3\sigma$ bounds to determine statistical significance (as a function of the number of available pairs having separation $D$) are also shown on the plot. No spatial correlation is seen, as all data points fall within the significance bounds. Though this plot shows the data from only one device size on one chip, similar plots for other device sizes and chips reveal the same conclusion.

$$\rho_{DVT,DVT\pm D}(D) = \frac{Cov(DVT, DVT \pm D)}{\sigma_{DVT}^2} \qquad (3.9)$$

The lack of spatial correlation indicates that $V_T$ variation is random, and we find that it is normally distributed. Figure 3-10 indicates that over 99% of the data points in this data set are indeed normally distributed. As supported by the measurement results at more scaled nodes, such as in [54, 55, 58], this leads to the conclusion that $V_T$ variation must be dominated by Random Dopant Fluctuation (RDF), even at the $0.18 \mu m$ technology node. Although Line-Edge Roughness (LER) should be considered as a possible cause of random variation, the $0.18 \mu m$ technology node is likely not affected by sub-resolution patterning effects to the same degree as a 90 or

Figure 3-10: $V_T$ distribution and normal probability plot for a single die

$65nm$ process. Furthermore, by measuring sub-threshold currents and reducing $V_{DS}$, we reduce the DIBL effect and minimize any current variations due to channel length variation. Additionally, oxide thickness is relatively well controlled and variation in this parameter is likely to be more spatially smooth. We note that smaller device sizes, including any excluded by our ADC resolution screening procedure, would be expected to be even more susceptible to RDF; an interesting result here is that RDF is discernable and dominates even for larger device sizes.

## 3.3.5 Die-to-Die Correlation

Given the lack of any significant intra-die spatial correlation, inter-die correlation is only expected if the standard deviation of the mean shift between each die is considerably larger than the within-die standard deviation (i.e., $\sigma_{die-to-die} > \sigma_{intra-die}$). Computing this requires the absolute $V_T$ for the reference device on each chip, which necessitates current measurements to determine where the $I_D$ vs $V_{GS}$ curve breaks from a straight line on a semilog plot. However, due to the limited dynamic range of the on-chip ADC, we were unable to measure currents significantly larger than $600nA$, making it impossible to determine an absolute $V_T$ for each reference device. Nevertheless, we are still able to compute correlations in the spatial variation patterns between arbitrary pairs of die. Once again, we choose a single device size to analyze

and compute the correlation coefficient, $\rho(i,j)$, between pairs of die using Eq. 3.10, where $DVT_{i_{(x,y)}}$ and $DVT_{j_{(x,y)}}$ are the $\triangle V_T$ of identical devices located at $(x,y)$ on die $i$ and die $j$, respectively, and $\mu_i$, $\mu_j$, $\sigma_i$, $\sigma_j$ are the means and standard deviations of $\triangle V_T$ of all devices of the given size on dies $i$ and $j$, respectively. Figure 3-9(b) shows the correlation coefficient as a function of all $\binom{36}{2} = 630$ pairwise combinations. All data points fall within the $3\sigma$ significance bounds indicating no significant variation pattern similarity between pairs of dies. Larger device sizes show the same results.

$$\rho(i,j) = \frac{Cov(i,j)}{\sigma_i \sigma_j} = \frac{\frac{1}{N}\sum_x \sum_y (DVT_{i_{(x,y)}} - \mu_i)(DVT_{j_{(x,y)}} - \mu_j)}{\sigma_i \sigma_j} \tag{3.10}$$

The previous two subsections have shown that within-die $V_T$ variation is both spatially uncorrelated and random (no repeatable within-die pattern) from die-to-die, implying that knowing the variation between two devices gives no further information about two similarly spaced devices on the same chip or even the same two devices on a different chip. We can now also update Figure 2-10 to include $V_T$ variation at or near the origin of the axes, as shown in Figure 3-11.



Figure 3-11: Decomposition of within-die spatial variation, including $V_T$ variation.

# 3.4 Design and Modeling Implications

These results have powerful implications for both circuit designers and the modeling community. In order to determine the importance of the lack of correlation in $V_T$ variation, we must determine the sensitivity of a given circuit or design to $V_T$ variation. For example, it is well known that sub-threshold circuit designs are highly susceptible to $V_T$ variation, while channel length variation has minor impact in comparison. In general, computing circuit sensitivity to individual process parameters is not easily done analytically. However, we can consider a simple inverter where the power-supply to the inverter is scaled from sub-threshold to well above $V_T$ to illustrate the differing circuit sensitivities.

## 3.4.1 Circuit Modeling

For a simple inverter operating above-threshold, the propagation delay for a falling transition can be modeled as in Eq. 3.11. However, when the inverter is operating sub-threshold, Eq. 3.11 is modified to Eq. 3.12.

$$\tau_{PHL} = \frac{C_L V_{DD}}{K_N \frac{W}{L} (V_{DD} - V_T)^\alpha} \tag{3.11}$$

$$\tau_{PHL} = \frac{C_L V_{DD}}{I_o e^{\frac{V_{DD} - V_T}{n V_{TH}}}} \tag{3.12}$$

where $1 \leq \alpha \leq 2$ for modern processes and $K_N$ is a constant determined by the process technology. Given the correlation coefficients for channel lengths ($\rho_{L1,L2}$) and threshold voltages ($\rho_{V_{T1},V_{T2}}$) between the NMOS transistors in two arbitrary inverters, it should be possible to compute the correlation coefficient for their propagation delays. In general, this is a difficult problem to solve analytically, but Monte Carlo simulations provide insight. To setup these simulations, we use the *delvto* parameter in Hspice to apply completely random variation in $V_{T_o}$ and we modify gate lengths in the spice deck using correlated random variation of $\rho_{L1,L2} = 0.9$. This was done in accordance with our results showing lack of correlation in $V_{T_o}$ and the spatial

74

(a) Correlation in inverter propagation delay

(b) Correlation in delay between RO and critical path of 64-bit KS adder

Figure 3-12: Circuit performance correlation as a function of $V_{DD}$

correlation results in $L$ from [70][4] to highlight the differing impacts of correlated versus uncorrelated variation in different circuit operating regimes.

We see in Figure 3-12(a) that the overall correlation contains some interesting characteristics. When $V_{DD} > 2V_T$, the correlation asymptotically grows toward $\rho_{L1,L2}$. However, around $V_{DD} = 2V_T$, the correlation falls precipitously, indicating that the gate overdrive is no longer sufficient to mask out RDF-dominated $V_T$ variation. Since Eq. 3.11 becomes less applicable when $V_{DD} \approx V_T$, the computed $\rho_{\tau_{PHL1},\tau_{PHL2}} = 0$, but below $V_T$ it settles to an extremely weak correlation, indicative of the dominant role of $V_T$ in determining sub-threshold current.

For simple logic gates, assuming channel length and $V_T$ variation are independent, one can safely assume that $\rho_{\tau_{PHL1},\tau_{PHL2}} \approx \rho_{L1,L2}$ when $V_{DD} >> V_T$. In contrast, in the sub-threshold region it is safe to assume $\rho_{\tau_{PHL1},\tau_{PHL2}} \approx 0$ due to the lack of correlation in $V_T$ variation.

This result demonstrates the importance of modeling correlations in each device parameter individually: the correlation of the particular circuit output metric is, in general, a non-linear function of the correlations in each device parameter. Fur-

---

[4]At the beginning of this chapter, we noted that within-die spatial correlation in channel length was an artifact of systematic variation. However, this decomposition, which yielded the aforementioned insight, had not yet been performed and only subsequently done in [16]. Nevertheless, the insights gained from this choice of $\rho_{L1,L2}$ can be applicable where systematic variation happens to be highly spatially correlated.

thermore, the overall circuit performance correlation is highly dependent on circuit operating region and the circuit's sensitivity to individual device parameters within the operating region. For regions of operation where correlation is low or insignificant, methods for modeling variation as uncorrelated IID statistics are appropriate, using Monte Carlo or distribution propagation approaches.

## 3.4.2 Circuit Design

Designers face many decisions in today's complex circuits. Increased variation means that designers must now consider how best to ensure robust circuit operation. One approach is to design the circuit to operate correctly given the range and characteristics of known or estimated device variations. A more aggressive approach is to consider active variation sensing and compensation strategies. The results of the previous section imply that the designer must carefully consider the operating region of the circuit when evaluating either robust design or active compensation strategies to counteract variation.

In the above-threshold region where channel length variation dominates circuit performance variation, assuming that channel length variation is spatially correlated, variation sensors can be used to detect variation. Active compensation, such as back-biasing circuits, can then be used to counteract the detected variation. Given high spatial correlation ($\rho > 0.8$) within a given radius, only one such variation sensor and compensation circuit is required within this radius.

In contrast, sub-threshold operation results in no significant spatial correlation due to $V_T$ or L variations. Additionally, Section 3.3.3 showed that the variance was not significantly related to distance. Consequently, a designer can only make use of the fact that the variance decreases with increasing device area, as increasing the size of a device effectively averages out variation due to Random Dopant Fluctuation. However, doing so results in negative power scaling due to increased total capacitance, and so there exists a trade-off between yield and power. In [109], the authors analyze the effect of increasing device width on the minimum energy point of sub-threshold operation, showing that up-sizing for a constant yield has a negative effect on the

minimum energy point. Another method of averaging out variation is to increase logic depth as in [110]. By increasing the number of devices in a logic path, the relative standard deviation of the propagation delay of the entire path decreases.[5] This technique can prove useful in designs where operating frequency is not the primary metric.

In Ultra-Dynamic Voltage Frequency Scaling (UDVFS) systems, such as those in low-power systems like mobile devices or sensor networks, both sub-threshold and above-threshold operation are used [111, 112]. Such systems often utilize replica critical paths or variation sensors to determine the appropriate frequency to operate at for a given power-supply voltage. To determine how correlated these replicas are to the actual path they are monitoring while scaling $V_{DD}$, we performed Monte Carlo simulations with correlated channel length variation of a 7-stage ring-oscillator (RO) and a critical path of a 64-bit Kogge-Stone (KS) adder. Figure 3-12(b) shows the results of a 1000-point Monte Carlo simulation with 65nm Predictive Technology Models [113], where $\rho_{Li,Lj} = 0.9$, $\rho_{V_{Ti},V_{Tj}} = 0$, $A_{V_T} = 5mV \cdot \mu m$, and $\sigma_L = 5\%$.

The overall correlation in delay between the two paths is $\approx 0.95$ for $0.7V \leq V_{DD} \leq 1.2V$. The reason this is greater than $\rho_{Li,Lj}$ is that the speed of one stage is directly dependent on the load provided by the following stage, which, in addition to the correlated variation in channel length, increases the overall correlation. In the region of $0.3V \leq V_{DD} \leq 0.7V$, the overall correlation between monitor and circuit decreases quickly, indicating that performance of the monitor/replica is no longer indicative of the performance of the critical path and should not be used to determine correct operating frequency unless a large guard-band is applied. Such a guard-band would undoubtedly subtract from increased energy efficiencies achieved by moving to lower supply voltages. A more robust method of controlling operating frequency in such systems is to detect logic errors in potential critical paths and slow down the frequency until timing errors are no longer detected [91, 114].

---

[5]The absolute standard deviation increases by $\sqrt{n}$ as described in Section 2.3 but the mean increases by a factor of $n$, so the standard deviation relative to the mean ($\frac{\sigma}{\mu}$) decreases by a factor of $\sqrt{n}$.

# 3.5 Summary

This chapter motivated the need to characterize variation at the device level to afford circuit designers a means of modeling and simulating the effects of variation in transistor parameters. The need to characterize threshold voltage variation became apparent as it is, along with channel length, a parameter capable of significant impact to circuit performance, especially in low-power applications. Analysis of the significance of threshold voltage in the sub-threshold region of MOSFET operation led to the design of a test-chip capable of isolating intrinsic threshold voltage variation by measuring these leakage currents.

With a designed capability to measure thousands of transistors in a dense array, the test-chip allowed characterization and spatial decomposition of threshold voltage variation. The measured data shows that intrinsic threshold voltage variation is, with high-likelihood, a truly random process that can be described accurately by using a Gaussian distribution with parameters based on transistor size. Furthermore, spatial decomposition and analysis of the data showed no within-die spatial correlation and no systematic, repeatable within-die variation pattern.

These results, in combination with channel length variation, were used to determine the effect on common digital circuits. The resultant Monte Carlo simulations showed that the effect on circuit performance, and appropriate variation mitigation schemes, is highly dependent on the application (e.g., high-performance versus low-power). This motivates correlating these simulation results with actual circuit performance variation and is the subject of the next chapter.

# Chapter 4

# Digital Circuit Performance Variation

Digital circuit design employs abstraction to a great deal. It is fitting then to begin to abstract away the variation in individual transistors and their parameters and think about variation in the context of digital logic gates or even entire circuits. This chapter begins by motivating the need for such abstraction (Section 4.1) and then, in Section 4.2, explores a test-chip architecture and relevant circuits capable of high-resolution (both temporal and spatial) variation characterization at the digital circuit level. The data analysis found in Section 4.3 provides key insights into the spatial decomposition of variation at the circuit level. In particular, the data again show no within-die spatial correlation, but at this level of abstraction we find systematic components of variation.[1] We also compare these results to the results predicted in the last section of the last chapter when combining both $V_T$ and $\triangle L$ variation to predict correlation in digital circuit perforance variation. Lastly, Section 4.4 discusses how the measured data necessitate unique variation mitigation strategies based on the operating mode of the digital circuits.

---

[1]Recall in the last chapter that there was no correlation in the $V_T$ variation pattern from one die to another, indicating no systematic component of variation.

# 4.1 Motivation

In the ideal design flow, digital designers can code desired circuit behavior using a high-level hardware descriptor language and then achieve implementation of the entire physical design using completely automated synthesis, placing and routing tools. The timing models used in these tools utilize abstraction as well; rather than modeling the behavior of individual transistors in each simulation, simpler gate-level models are utilized. Although not as accurate as simulating individual transistors, modern gate-level models can yield results within 5% of more detailed transistor-level models, while providing great enhancements in speed and productivity.

Enabling modeling of variation in digital circuit perfomance means providing variation models at the same level of abstraction. Variation models that can accurately capture variation at the gate-level can be directly plugged into the relevant tools. Going a step further, characterization of common digital circuits provides designers the intuition necessary to guide robust design as well as concrete physical data capable of verifying the results of automated tools.

As decomposition of the spatial components of variation increases in significance, understanding circuit performance correlation enables first, more accurate modeling, as evidenced by the emergence of statistical timing models which incorporate spatial correlation [98, 105, 99] despite not having manufacturing data to validate such models, and second, design of appropriate mitigation techniques. Specifically, if nearby circuits are strongly correlated ($|\rho| \approx 1$), techniques involving "replica" circuits can be employed as both monitor circuits ("canaries") as well as within feedback loops capable of active compensation. However, when circuits are weakly or not at all correlated ($|\rho| \approx 0$), "in-situ" techniques are likely necessary, as replica circuits may not accurately mimic the performance of the actual circuits under consideration. Furthermore, correlation data, and in particular any spatial dependencies, enable determination of spacing criteria between monitor circuits and potential critical paths in the design to ensure high correlation.

The next section describes the architecture of an all-digital test-chip capable of

providing spatial variation data at the circuit level.

## 4.2   Test Circuits & Chip Architecture

The variation test-chip is composed of a large number of replicated blocks containing various digital circuits and all-digital measurement circuitry. The following subsections further detail the architecture of each component.

### 4.2.1   High Frequency Test Circuits

To extract spatial correlation of high-frequency circuits, we array 80 nominally identical "adder-blocks" arranged in a 9x9 matrix as shown in the die photo of Figure 4-1(a) and occupying a $4mm^2$ area. Each of these blocks contains common circuits found in modern product designs, namely a 64-bit adder, canary ring oscillators of various types and frequency measurement circuits as shown in Figure 4-1(b). Details of each of these components are provided below.



(a) Die photo of test-chip implemented on IBM's $90nm$ CMOS process.

(b) Adder-block floorplan showing all circuits within the replicated block.

Figure 4-1: Arrayed Kogge-Stone adders instrumented for internal delay measurement.

## Oscillating 64-bit Kogge-Stone Adder

The first circuit included in each "adder-block" is a 64-bit Kogge-Stone adder. Kogge-Stone adders belong to a class of adders known as parallel-prefix architectures which enable fast, efficient addition by pre-calculating carry signals. The critical path in such adder structures is logarithmically related to the number of bits being added rather than linear in the worst-case, thus making wide adders feasible. Such adder structures are common building blocks in the datapaths of modern ASICs and microprocessors, making them good candidates for spatial variation analysis.

To enable simple and efficient measurement of maximum operating speed, the adder inputs are configured such that the outputs oscillate, requiring only a very simple frequency counter for adequate characterization of performance. This is accomplished by setting one of the inputs to all ones ($A\langle 63 : 0\rangle = 1$), the other input to all zeros ($B\langle 63 : 0\rangle = 0$), and connecting an inverted version of the carry-out signal to the carry-in input ($C_{in} = \overline{C_{out}}$), as shown in Figure 4-2.



Figure 4-2: Oscillating 64-bit Kogge-Stone adder

All gates in the design are a custom standard cell implementation, each sized (using the Hspice optimizer) such that they meet an output rise/fall time constraint ($50ps$ in this implementation) for a given multiple of the minimum size load (i.e., a 1X cell must have a $50ps$ output rise/fall time for a load of $5fF$, for a 2X cell the load is increased to $10fF$, etc.). Carry-propagate cells (AND gates) and carry-generate cells (AND-OR gates) constitute the largest part of the adder, with the remainder

of the gates being primarily XOR gates to perform the addition. The standard cell implementation allows for ease of layout due to the fixed pitch and the ability to abut cells without design rule violations. Although layout and routing were both manually done, this is representative of modern design flows for high performance blocks as custom layout and routing often results in higher performance than completely automated flows.

## Canary Ring Oscillators

To capture correlation between different circuit structures, we include 16 ring oscillators (ROs) around each adder, since they are often used as monitor circuits due to their simplicity and small area overhead. The ROs used are: INV{9,11,13,15}, NAND{9,11,13,15} and NOR{9,11,13,15}, denoting the type of gate and number of gate delays. The INV9, NAND11, NOR13 and INV15 are duplicated for a total of 16 ROs per adder. The NAND and NOR gates are two-input gates with inputs shorted to produce an inverting gate. These differ from the AND-OR gate used in the KS adder, which allows for quantification of correlation between disparate gate types and transistor stacks. The 16 ROs are divided into four blocks and interspersed between other digital logic, as shown in Figure 4-1(b), to ensure that these "canary" circuits are within the context of circuits normally found in modern digital designs.

## Frequency Measurement Circuits

Two 32-bit asynchronous frequency counters, made up of simple toggle flip-flops (Figure 4-3), measure the oscillation frequency of the adder and ROs. An additional two frequency counters are also included for use in measuring the delay of each individual adder bit relative to the first bit in the adder, detailed in the following section. These frequency counters provide a simple solution to high-resolution digital measurement, allowing completely digital read-out of all relevant frequencies and delays. The outputs of the frequency counters are multiplexed onto a single bus spanning each row of the chip, and these row buses are multiplexed onto the chip outputs.

Figure 4-3: 32-bit asynchronous frequency counter comprised of toggle flip-flops.

## 4.2.2 High-Resolution Digital Delay Measurement

Capturing variation data with higher spatial and temporal resolution than possible with simple ROs requires alternative techniques. Specifically, we seek a highly-scalable, all-digital measurement technique capable of sub-picosecond delay resolution occupying small area. In the following sections we describe the theory and operation of a random-sampling technique, in the context of measuring delays between bits of the adder, meeting these criteria.

**Random Sampling**

The critical path in the adder is from $C_{in}$ to $C_{out}$, resulting in all bits oscillating at an identical frequency but with delays determined by the logarithmic structure of the KS adder, giving 64 out-of-phase oscillator taps. Quantifying this phase-delay is equivalent to quantifying the difference in delay between each bit, and allows for variation analysis with improved spatial resolution. Delay measurement is done by randomly sampling the signals and counting the number of occurrences when one of the signals is logic high and the other simultaneously low, or vice-versa, dividing by the total number of samples taken and multiplying by the signal period as shown in Eq. 4.1.

$$D = \tau_{period} \frac{N_{up} - N_{dn}}{N_{tot}} \tag{4.1}$$

Uses of random-sampling techniques for this purpose are not new [115, 116, 117]. However, in each of those implementations, an XOR gate is used to determine when one of the signals is logic high and the other logic low. XOR gates can limit the minimum detectable delay since some minimum delay is necessary to register a clear

logic level at the output. In our implementation, we transform the measurement to be that of the difference in two pulse widths using a Phase-Frequency Detector (PFD), discussed in the following section, which eliminates this constraint.

For this technique to be accurate, the random samples must be uniformly distributed across all points in the sampled signal cycle as non-uniformity in this distribution results in non-linearities in the measured delays. Mathematically, the operation shown in Figure 4-4(a) can be modeled as:

$$S_{N+1} = S_N + X_N = \sum_{i=1}^{N} X_i \qquad (4.2)$$

where $S_N$ is the timing of the Nth sample, $X_i$ is the time between samples and is a random variable, and $t_{per}$ is the sampled signal period. We also define $P_N = mod(S_N, t_{per})$ as the position of the Nth sample in the sampled signal cycle. If the average sampling period ($\mu_X$) and the sampled signal period ($t_{per}$) are co-prime, it is clear that the distribution of $P_N$ will be uniform even without random edges. However, if they are not co-prime, Eq. 4.2 describes a random walk, which we simulated in Matlab using a distribution of random periods taken from Hspice simulations of the LFSR-controlled RO. The results of this simulation, in Figure 4-4(b), empirically show convergence to a uniform distribution as $N \to \infty$. With $N > 10^8$ samples the non-uniformity of the distribution, and thus non-linearity in the measured delay values, is sufficiently small.

The framework in [117] can be used to find standard errors and confidence intervals on the time resolution of this technique given some number of samples. Conversely, it can also be used to compute the number of samples required for a desired resolution given desired confidence intervals. Assuming equal probability of a sampling instant falling anywhere in a clock cycle, Eq. 4.3 shows how the number of samples required, $N$, depends on $z_c = \sqrt{2}erf^{-1}(CI)$ which gives the area underneath an appropriate Gaussian curve for a certain confidence interval, $CI$, and $p$ and $P$, which are the

(a) Random sampling to determine pulse widths (b) Convergence of sampling instant to a uniform distribution within a clock cycle

Figure 4-4: Graphical depiction of random sampling

actual and observed duty cycles, respectively.

$$N = \left(\frac{z_c}{p - P}\right)^2 P(1 - P) \tag{4.3}$$

As shown in Figure 4-4(b), $N > 10^8$ samples will satisfy the uniform distribution condition that satisfies the assumption of this analysis. Section 4.3.3 will show that the number of samples taken in our random sampling implementation is $\approx 1.7 \times 10^8$, satisfying the constraint. It should also be noted that the time resolution of this technique can be arbitrarily increased with large enough $N$, which would result in smaller errors between actual and observed duty cycle $(p - P)$.

### Measurement Circuits

Most high-resolution, sub-picosecond measurement techniques require complicated circuitry capable of exquisite timing. The nature of random sampling techniques, in making use of a large number of samples spread over many cycles, relaxes the measurement circuit constraints substantially. In this implementation there are two major circuits: a Phase-Frequency Detector (PFD) which converts a delay difference into a pulse-width difference and a random sampling clock generating circuit.

Removing the constraint on minimum measurable delay is achieved by transform-

ing the measurement to that of the difference in two pulse widths. Since we are attempting to measure the delay difference between bits of the KS adders and the adders are in an oscillating configuration, the oscillating frequency and pulse widths remain constant for the duration of operation. This situation lends itself nicely to sampling over a large number of cycles as the inherent properties of the sampled signal remain constant, or stationary, over time.

A phase-frequency detector as shown in Figure 4-5(a), normally used in phase-locked loops, is used here to convert between a delay difference and a pulse-width difference. The circuit is comprised of two modified True Single-Phase Clock (TSPC) flip-flops and a self-resetting circuit. For simplicity in explanation, we assume that the signal that arrives first, or the early signal, is $Sum <x>$, the input to the top TSPC flop; however, the following explanation, graphically depicted in Figure 4-5(b), can be completely reversed with no loss in circuit functionality. When the early signal arrives it triggers a rising edge on the $UP$ signal. Similarly, when the late signal arrives at some point later, a rising edge is triggered on the $DN$ signal. The arrival of the late signal simultaneously triggers the self-resetting circuitry connected to the outputs of both flops, resulting in *simultaneous* falling edges on both the $UP$ and $DN$ signals.[2] As the timing diagram shows in Figure 4-5(b) by virtue of the simultaneous falling edge, the circuit has transformed a delay measurement into a pulse-width difference measurement. Moreover, the difference of these two pulse widths can be arbitrarily close to zero, eliminating any bound on minimum measurable delay.

Since the $UP$ and $DN$ pulses have differing pulse-widths, it is these signals that are randomly sampled using a simple D Flip-Flop, clocked by a random clock. Generating this random clock is done by attaching the outputs of a Linear Feedback Shift Register (LFSR) to the tri-state controls of inverters in a 5-stage ring-oscillator as shown in Figure 4-6. Each stage of the ring-oscillator is variable drive-strength; since the LFSR produces pseudo-random bit sequences, the tri-state controls of the ring-oscillator are enabled in a pseudo-random fashion, in effect generating the random clock we require.

---

[2]Variation and mismatch between the flops result in a non-simultaneous reset and a corresponding static offset. However, with the measurement technique discussed in Section 4.3.3, this offset can be accounted for and eliminated.

(a) PFD [118]                    (b) Timing diagram

Figure 4-5: PFD and associated timing diagram

Since LFSRs can repeat the pseudo-random pattern if not designed with enough bits, this implementation is comprised of 63 bits, which theoretically is sufficient to avoid repetition over many years of operation at the designed operating frequencies.



Figure 4-6: Random clock generation using a LFSR and variable drive-strength ring-oscillator.

Hspice simulations show the distribution of clock periods generated by this circuit (Figure 4-7). The generated clock periods are indeed randomized with a distribution that can be closely approximated by a Gaussian distribution except in the tails. Nevertheless, we are more interested in the distribution of sampling instants within a clock cycle, which the previous section showed to be uniform. Despite the relatively

88

simple circuitry and the fact that the clock periods form an approximate Gaussian distribution, with enough samples, the "drifting" nature of these periods results in a uniform distribution of the sampling instants within a clock cycle.



Figure 4-7: Distribution of clock periods of random sampling clock

A single PFD is included within each replicated "adder-block" and shared between all bits of the adder. The first sum bit of the adder is directly connected to one of the inputs of the PFD and all 64 bits of the adder are mux-able into the other input of the PFD. This setup allows subtraction of any offsets (e.g., due to multiplexing or variation in the two halves of the PFD) or other non-idealities by measuring the first sum bit versus the directly connected version of itself and subtracting that result from all other measurements.

### 4.2.3 Body-Biasing Circuits

The last set of circuits included as part of this test-chip are intended for testing adaptive body-biasing as a variation mitigation technique. Although such a study was performed in [78], it is well known that the efficacy of adaptive body-biasing should decrease in scaled technologies. This is a result of the $\gamma$ coefficient of the last term in Eq. 3.1 being directly proportional to the oxide thickness as shown in Eq 4.4 (recall that $C_{ox}$ is inversely proportional to the oxide thickness from Eq. 2.5 in Section 2.2.2), whereas there is a square root dependence of $\gamma$ on doping concentration. To maintain gate control of the channel as doping concentrations increase, the oxide

thickness is simultaneously decreased, typically resulting in a smaller $\gamma$.

$$\gamma = \frac{\sqrt{2q\epsilon_s N_{ch}}}{C_{ox}} \qquad (4.4)$$

The ability to adaptively body-bias both PMOS and NMOS transistors of the KS adders is added by using IBM's triple-well process which allows NMOS transistors to be isolated from the substrate. Large pass-gate transistors, shown on the left edge of Figure 4-1(b), are used as body-bias multiplexers, allowing selection of one of eight body-bias voltages generated external to the chip. Simple flip-flops control the multiplexers and are set using software-controlled signals.

## 4.3 Data Analysis

The test-chip described above was fabricated in IBM's 90nm triple-well process technology, as pictured in Figure 4-1(a), and frequency measurements were carried out on 41 chips from two wafers. This section describes the measurements and quantifies their spatial decomposition.

### 4.3.1 Variation Measurements

Variation measurements show that the effect of variation is larger at lower power-supply voltages (Figure 4-8). As $V_{DD}$ is decreased below 0.7V, variation increases dramatically. Moreover, the within-die component becomes a larger and larger component of the overall variation. While noise is typically a concern at low voltages, all of the measurements are taken by allowing the circuits to run freely for $> 10^8$ cycles, averaging out any (assumed white) noise.

As expected, ROs with smaller gates (e.g., INV) and fewer number of delay stages show greater within-die variation (due to less averaging, see below) than ROs with either larger gates (e.g., NOR) or more delay stages. At the extremes, the INV9 RO is most variable ($\sigma_{WID} = 4.9\%$ at 0.5V and $\sigma_{WID} = 1.13\%$ at 1.2V) and the adder, being the largest circuit, is least variable ($\sigma_{WID} = 2.8\%$ at 0.5V and $\sigma_{WID} = 0.82\%$

at 1.2V). Die-to-die and total variance show the same general trends except that the NAND ROs are more variable than the INV ROs — we believe this may be due to larger NMOS than PMOS global variability as NAND gate delays are typically dominated by the series NMOS stack.



Figure 4-8: Variation as a function of $V_{DD}$

Decomposing the variation into die-to-die and within-die components reveals that not only does the fraction of variation attributable to within-die variation increase with decreasing voltage, but perhaps more importantly, for some circuits, namely the INV ROs, the within-die and die-to-die components are roughly equal at those low voltages. Even in the "best" cases, the within-die component increases from 25% of the die-to-die component when $V_{DD} > 1V$ to 50% at low voltages ($V_{DD} = 0.5V$). The asymmetric increase in variation components significantly impacts the design of variation mitigation schemes, especially those operating in low-power and low operating voltage regimes or dynamically scaling systems, as will be detailed in Section 4.4.1.

## 4.3.2   Spatial Variation & Correlation

The large number of replicated circuits within each die allows extraction of spatial correlation. Within-die spatial correlation is computed in the same manner as in Chapter 3.3.4, except that in this case the mean of each die is subtracted and all

circuits separated by a distance, $d$, over all die are included in the computation. This increases the number of data points per distance, thereby increasing the statistical significance of the computed results. At the closest separation distance of $170\mu m$ between adjacent adders, there are nearly 1476 (41 chips $\times$ 36 adjacent adders) independent pairs of adders that contribute to the correlation computation, while at the farthest separation distance of $2mm$ there are only 82 (41 chips $\times$ 2 pairs) independent pairs of adders.



Figure 4-9: Adder spatial correlation as a function of $V_{DD}$

We find that even at high $V_{DD}$, there is no discernible *within-die* spatial correlation (Figure 4-9). The dip at the $V_{DD} = 1.2V$ and $2mm$ separation distance corner is likely due to limited data, as there are only 82 pairs of adders at this separation distance over 41 die. Furthermore, all points are within a 95% confidence interval, giving additional credence to this argument. Lack of within-die spatial correlation does not imply a lack of systematic within-die variation: rather, it simply means such variation is not a function of separation distance. Instead, this implies that the statistics of adder frequency (e.g., mean and standard deviation relative to chip-mean) is dependent on position within the die, as evidenced in Figure 4-10. Indeed, a systematic within-die pattern in adder frequency is noticed in Figure 4-10(a). To elucidate this further, consider equidistant pairs of adders $\{(6,6), (6,7)\}$ and $\{(7,6), (7,7)\}$ that exhibit both unequal and directionally opposite frequency differences, $\triangle F(\{6,6\}, \{6,7\})$ and

$\triangle F(\{7,6\}, \{7,7\})$, which implys lack of correlation.



(a) Average within-die spatial variation pattern (Per location $\triangle F$ from die mean averaged over all die).

(b) Location-dependent standard deviation of frequency difference from chip-mean $\left(\frac{\sigma}{\mu}\right)_{\triangle F}$ over all die.

Figure 4-10: Systematic within-die variation

Since this systematic pattern introduces position-dependent variation statistics, the pattern should be removed prior to any spatial correlation computation to ensure that the residue is stationary. This is necessary, as computing correlation implicitly assumes that all samples are from the same probability distribution.[3] When the systematic pattern shown in Figure 4-10(a) is subtracted (after subtracting out die means), the newly computed spatial correlation is shown in Figure 4-11. Comparing the two figures, there is little difference owing to the mostly random nature of the systematic pattern. The small peaks that were present due to the slow bottom row of adders have now been flattened out, showing conclusively that there is no spatial correlation in within-die circuit variation.

The distinction between position-dependent systematic variation and separation distance-dependent variation or spatial correlation again allows us to update the axes we introduced in Section 2.3. As mentioned above, digital circuit performance variation has a systematic component but no spatial correlation component and thus belongs somewhere directly on the x-axis as shown on Figure 4-12.

Nevertheless, Figure 4-9 also shows strong die-to-die correlation at high $V_{DD}$ but

---

[3]This was not necessary for the spatial correlation computations in Chapter 3 as there was no correlation in the variation patterns between die reflecting no, or very weak, systematic variation.

Figure 4-11: Adder spatial correlation as a function of $V_{DD}$ with systematic components removed.

decreasing with lower $V_{DD}$, indicating that the effect of random variation increases at lower power-supply voltages. This is consistent with $\sigma_{D2D} \gg \sigma_{WID}$ at higher voltages, but the relative fraction decreasing at lower voltages as seen in Figure 4-8. Since the within-die variation is random (spatially uncorrelated and even the systematic component not containing any apparent order), die-to-die correlation should decrease as the within-die component increases in relative strength. These results are also consistent with the effect of threshold voltage ($V_T$) variation, which is dominated by Random Dopant Fluctuation (RDF), increasing as gate overdrive decreases [119]. Furthermore, die-to-die correlation shows only weak dependence on separation distance. Qualitatively, this means that when die-to-die correlation is strong, knowledge of how two identical circuits differ from one die to another enables strong inference when comparing any other circuits (of the same type) between those die, regardless of how far those circuits may be separated from the original circuit. However, if voltages (and thus correlation) are decreased, such inferences become increasingly weak.

While Figure 4-9 only shows adder correlations, the spatial correlation results for all ring oscillators are similar, with the notable exception that die-to-die correlation decreases with decreasing $V_{DD}$ more quickly with smaller circuit size. As an example, the INV9 based RO has a correlation coefficient, $\rho$, of 0.55 at $V_{DD} = 0.5V$ compared

94

Figure 4-12: Decomposition of within-die spatial variation, including digital circuit performance variation.

with $\rho = 0.65$ for the INV15 (shown in Figure 4-13) and $\rho = 0.75$ for the adder, consistent with smaller circuits being more susceptible to random variation sources as less averaging occurs. On a die-to-die scale, these results closely match the results of the simulations from the end of the last chapter (Section 3.4.1) which showed that correlated channel length variation could result in strong correlation at high $V_{DD}$ but degrades rapidly as $V_{DD}$ is scaled downward to save power.



(a) 9-stage inverter RO



(b) 15-stage inverter RO

Figure 4-13: Inverter spatial correlation plots showing different decreases in die-to-die correlation with decreasing voltage based on number of stages.

Since ring oscillators are commonly used "canary" devices, thought to predict the performance of circuits situated nearby, analyzing the cross-circuit correlation is im-

portant as well. Extrapolating from the earlier results showing weak or no within-die spatial correlation, a lack of cross-circuit correlation is expected. Figure 4-14 contains scatter matrices for each pair-wise combination of circuit types. Perfect cross-circuit correlation ($\rho = 1$) would be manifested by a tight distribution, along a positive diagonal line ($y = x$), of all the scattered points within that individual axes. Similarly, perfect anti-correlation would be manifested in the opposite manner, i.e., a distribution along a negative diagonal line ($y = -x$). Lastly, circular scatterings reflect no correlation at all, which is the case in all of the pair-wise circuit combinations. Scatterings that form tighter distributions along either the horizontal or vertical axes indicate only that there is less frequency variation in one of the circuits in the pairing (as evidenced by the accompanying histograms), not any correlation. These results are for high voltage ($V_{DD} = 1V$) and the closest separation distance possible for each circuit pairing; each circuit is compared to the circuits in the same "adder-block." Repeating this calculation for larger separation distances or other voltages reveals identical results, as is expected from the earlier "same-circuit" spatial correlation results.

### 4.3.3 Adder Bit Delays

Using the all-digital delay measurement circuits described in Section 4.2.2, we measure the delay of each bit within each KS adder, relative to $Sum\langle 0 \rangle$ of the same adder. All 64 bits, including $Sum\langle 0 \rangle$, are muxed into a single PFD. By using the same PFD for all bits, offsets and other non-idealities due to mismatch in the PFD structure can be measured (by measuring the delay between the muxed and non-muxed versions of $Sum\langle 0 \rangle$) and subtracted from subsequent measurements.

Measurements over all 80 adders on 40 chips are shown in Figure 4-15 with a post-layout extracted simulation for comparison. There is good agreement between measured delays and simulation: the upward shift between measured data and simulation indicates a slower process than nominal simulation, and is consistent with 20% slower frequency measurements than nominal post-layout simulations. The measured delay pattern is consistent with simulation for all but three of the bits, $Sum\langle 16, 32, 48 \rangle$,

Figure 4-14: Scatter matrices of each pair-wise combination of circuit types showing no cross-circuit correlation at $V_{DD} = 1V$. All circuit pairings are within the same "adder-block."

which are typically the fastest due to the logarithmic structure of a KS adder. However, in this layout, these three bits contain longer wires than the other bits, corresponding to the peaks in the post-layout simulation. While these three bits do have larger measured delays than nearly all of the other bits, the difference is not as large as in simulation, possibly due to slower transistors but "faster" wires, which would decrease the delay peaks formed by these three bits.



Figure 4-15: Bit delay measurements relative to $Sum\langle 0 \rangle$

Each data point in Figure 4-15 consists of no less than $1.7 \times 10^8$ random samples. Using the theoretical analysis in [117], we compute a possible observed error of approximately $200fs$, with a confidence level of 99.9999%. Combined with the general agreement between measured data and simulation, this computation gives high confidence in sub-picosecond accuracy of the measured data.

Due to the correlated structure of the adder, care must be taken when doing spatial correlation analysis using the measured bitslice delays. Since the critical path of the KS adder involves all bits, intuitively all bits will be correlated with each other to some degree. The logarithmic nature (in particular, radix 2, or log-base-2) of the adder implies that bits two away, four away, eight away, etc. from each other will be more correlated than other combinations of bits. Figure 4-16 shows the cross-correlation of bit-slices within the same adder and indeed reveals this logarithmic

Figure 4-16: Within-adder bit-slice cross-correlations revealing the logarithmic structure of the adder.

structure. The prominent, darker diagonals starting on the x-axis at bits 9, 17, 25, 33, 41 and 49 show stronger correlations between bit-slices 32, 16 and to a lesser degree 8 and 4 away from each other. The 32 bit-slice separation is most strongly correlated, as it is an input nearer the end of the higher bit-slice's path and arrives later than outputs of bits nearby, therefore more strongly influencing the delay of that bit-slice. Also noticeable, and expected, are strong cross-correlations in bit-slices 1-16 and 33-48 by virtue of the carry signal connections to logarithmically pre-compute the overall carry-out signal.

Since the adder circuit structure largely determines the correlation of bit-slices, within-adder correlation does not reveal much about spatial correlation due to process variation. Nevertheless, looking at bit-slices relatively close to each other but not as strongly connected, for example bit-slices 16 and 20, reveals much weaker correlation: $\rho = 0.4$ compared to $\rho > 0.8$ for bit-slices separated by 32. Such results show that the

relative significance of circuit structure and connectivity is far greater than process variation, and is likely the only significant source of spatial correlation.

Knowing that circuit structure and connectivity significantly impact the apparent correlation, the correlation between each bit-slice and its own adder performance should be strong. However, Figure 4-17 shows that each bit-slice is only moderately anti-correlated with the performance of the adder it partly constitutes. The anti-correlation is simply due to comparing delays with frequencies and the inverse relationship between the two. The moderate, rather than the expected strong, correlation is likely an artifact of the comparison being made: the phase delays are all relative to $Sum\langle 0 \rangle$ while the adder frequencies are absolute. Unfortunately, since we are not able to measure the absolute delay of the $Sum\langle 0 \rangle$ bit-slice, we can only speculate that with absolute delays a stronger correlation would be extracted.



Figure 4-17: Correlation between bit-slice delays and frequencies of the adder they are a part of. Anti-correlation is due to comparing delays versus frequencies — an inverse relationship.

Further spatial correlation analysis of bit-slices across different adders is also possible and, as expected from all previous results, shows no significant spatial correlation. Figure 4-18 shows the within-die correlation between two bit-slices picked at random, revealing strong correlation only at the distance corresponding to the two bits being part of the same adder.[4] At larger separation distances, there appears to be slight

---

[4]Although there are smaller distances on the plot, these smaller distances correspond to these bit-slices being part of different adders due to the layout of the adder.

anti-correlation. However, as we noted previously with respect to the adder frequencies, there are far fewer independent pairs separated by these distances; the stronger anti-correlation here is likely due to random chance with smaller numbers of samples, as evidenced by the much larger confidence intervals. Despite showing the cross correlation for only two bit-slices here, any arbitrary combination of bits reveals similar results.



Figure 4-18: Spatial correlation between bit-slices 19 and 49 (blue dots) and associated confidence intervals (triangles). Strong correlation is only noticed for bit-slices within the same adder.

## 4.4 Variation Mitigation

Lack of spatial correlation indicates that the magnitude of the absolute within-die variance and any systematic variation components should guide design of within-die mitigation schemes. At the die level, only the magnitude of the die-to-die variation component is of significance. Furthermore, the data suggests that variation mitigation strategies must be a function of the voltage/power domain in which circuits are operated. This dependence is explored in the next subsection, prior to comparing the efficacy of adaptive body biasing and adaptive power-supply voltage scaling for within-die, circuit-level variation mitigation.

## 4.4.1 Dependence on Performance Domain

In high-performance domains where gate over-drive is sufficiently large $(> 2V_T)$, although within-die variation is random and uncorrelated spatially, the absolute variation is also small: the standard deviation is less than 2% of the mean. Mitigation schemes involving small "monitor" circuits, such as ring oscillators or replica critical paths, require little margin to be effective. In the case of larger circuits, such as adders, a 2% margin is sufficient to have a 95% confidence level (corresponding to $2\sigma$) that the critical path will track the monitor or replica circuit sufficiently well.

In low-performance, low-power domains in which $V_{DD} \approx V_T$, variation as a percentage of the mean is significantly increased, by as much as $5\times$ relative to variation in high-performance voltage domains. Furthermore, within-die variation is a significant component of the overall variation seen. In combination, these two factors necessitate in-situ circuits capable of accurate measurement of timing data as the basis of a robust mitigation strategy. Without in-situ measurement capability, large margins are required for sufficient confidence intervals. For example, a replica of the adder critical path would require 6-10% margin for 95-99% confidence. Shorter paths or smaller circuits require that the margins increase to as much as 15-20%, making monitor or replica circuits infeasible.

Two common mitigation strategies adaptively change circuit voltages (the substrate or body voltage, or the power-supply voltage) to tune performance. Either, or both, of these techniques can be used in both high-performance and low-power domains as long as a domain appropriate measurement/sensing scheme is used. However, both have advantages and drawbacks that can be specific to the domain.

## 4.4.2 Adaptive Body-Biasing

Body-biasing modifies the threshold voltage of a device by adjusting the voltage on the substrate, or body, terminal of the transistor based on Eq. 2.6, where $V_{SB}$ is the source-to-body bias voltage. If the frequency of circuit operation is the inverse of the delay in Eq. 3.11 in the above threshold regime, then the sensitivity of circuit

frequency to body-biasing can be computed by taking the derivative of the inverse with respect to $V_{SB}$. Eq. 4.5 shows that the frequency sensitivity is only weakly dependent on $V_{SB}$, as $\alpha < 1.5$ for modern processes and will likely result in a close to linear frequency scaling with $V_{SB}$. However, as $V_{DD}$ scales lower toward $V_T$, Eq. 3.12 must be used instead and will show an exponential frequency sensitivity.

$$\frac{dF}{dV_{SB}} = \frac{-\alpha\gamma}{2\sqrt{2\phi_{Fp} + V_{SB}}} \cdot \frac{K_N\frac{W}{L}\left(V_{DD} - V_{FB} - 2\phi_{Fp} - \gamma\sqrt{2\phi_{Fp} + V_{SB}}\right)^{\alpha-1}}{C_L V_{DD}} \quad (4.5)$$

This test-chip includes the capability to adjust the body bias of both the PMOS and NMOS transistors of the KS adders due to IBM's triple-well process which allows for NMOS transistors to reside in separate wells.[5] Figure 4-19 shows how the performance (frequency) of the adder can be tuned by adjusting both PMOS and NMOS body-biases. In particular, there is a large difference in tunability based on performance domain: at $V_{DD} = 0.5V$, the frequency can be tuned by as much as $-65/+137\%$, but at $V_{DD} = 1.0V$ the tuning range is reduced to $-19/+27\%$. This is largely a result of increased gate-overdrive (e.g., the value of $V_{DD} - V_T$ is larger) in high-performance domains reducing the impact of changes to $V_T$, just as noted in the variation results presented in Section 4.3.1 and expected from the analysis above.

While body-biasing likely provides sufficient tunability to mitigate the impact of even large variation at either performance domain, the impact on leakage current and power may be the determining factor. As noted in Chapter 3, leakage currents are an exponential function of $V_T$ (Eq. 3.2) so adjusting $V_T$ using body-biasing can severely impact leakage power. Although there is no capability to measure individual adder leakage currents in our test chip, Figure 4-20 shows results of simulations of the adder: just as predicted, there is an exponential impact on leakage currents. Improving adder performance by 25% at $V_{DD} = 1.0V$, corresponding to $V_{SBP} = 0.5V$ and $V_{SBN} = -0.5V$, comes with the penalty of 200% additional leakage current as both PMOS and NMOS threshold voltages are lowered by the non-zero body-bias,

---

[5]More commonly, twin-well processes result in all NMOS transistors sharing the same substrate. Consequently, changing the substrate bias affects all NMOS transistors whereas a triple-well process allows for selective body-biasing of NMOS transistors.

(a) $V_{DD} = 0.5V$        (b) $V_{DD} = 1.0V$

Figure 4-19: KS adder frequency deviation from nominal (no body-bias) versus body-bias magnitude

increasing sub-threshold leakage. Similarly, in the low-power domain of $V_{DD} = 0.5V$, a 137% performance boost results in a 600% increase in leakage currents, again when $V_{SBP} = 0.5V$ and $V_{SBN} = -0.5V$.



(a) $V_{DD} = 0.5V$        (b) $V_{DD} = 1.0V$

Figure 4-20: Simulations of adder leakage current deviation from nominal (no body-bias) versus body-bias magnitude

## 4.4.3 Adaptive Power-Supply Voltage Scaling

Rather than attempting to adjust $V_T$ to tune performance, adjusting the power-supply voltage can potentially be more impactful with reduced leakage current overhead. To see why this is, we again differentiate the inverse of Eq. 2.3, but with respect to

$V_{DD}$ this time to get the result in Eq. 4.6. Again, the derivative seems to be weakly dependent on $V_{DD}$, which should result in a roughly linear scaling of frequency with $V_{DD}$. Since the leakage current is exponentially dependent on $V_{DS}$ and thus $V_{DD}$, the leakage current should scale exponentially, but with a smaller constant as the DIBL coefficient, $\eta$, is much smaller than 1 in Eq. 3.2.

$$\frac{dF}{dV_{DD}} = \frac{K_N \left(\frac{W}{L}\right) \left[(\alpha - 1) V_{DD} - V_T\right] (V_{DD} - V_T)^{\alpha-1}}{C_L V_{DD}^2} \tag{4.6}$$



(a) Mean measured adder frequency versus $V_{DD}$, normalized to $V_{DD} = 1.0V$    (b) Simulated adder leakage versus $V_{DD}$, normalized to $V_{DD} = 1.0V$

Figure 4-21: Adder performance as a function of power-supply voltage, $V_{DD}$

Indeed when the mean measured frequencies of the KS adders are plotted versus $V_{DD}$, there is a roughly linear relationship, as shown in Figure 4-21(a). Also seen in Figure 4-21(b) is the exponential scaling in leakage current as a function of $V_{DD}$ due to DIBL, as expected.[6] When comparing these plots to Figure 4-19 where the body-bias is scaled, we find that scaling $V_{DD}$ is more impactful than a similar scaling of body-bias to tune frequency. However, it is hard to ascertain which is more efficient in terms of leakage current/power. For an accurate comparison of the two methods, plots of relative change in leakage current versus relative change in frequency are used, as shown in Figure 4-22.

---

[6]The simulated data points are fit to a simple DIBL model ($I \propto e^{\frac{\eta V_{DD}}{nV_{th}}}$) and the fitted parameters are consistent with those given by the process.

(a) Body-biasing, nominal $V_{DD} = 1.0V$  (b) $V_{DD}$ scaling, nominal $V_{DD} = 1.0V$

Figure 4-22: Relative change in leakage as a function of relative change in frequency when body-biasing or $V_{DD}$ scaling are used for variation mitigation.

In Figure 4-22(a), adaptive body-biasing is used for performance tuning. Since both PMOS and NMOS transistors are biased, there are multiple leakage currents for each frequency setting, but only one combination results in a minimum leakage for each frequency. Nevertheless, when compared to power-supply scaling in Figure 4-22(b), it is evident that power-supply scaling can achieve larger frequency tuning range with less negative leakage current impact. When increasing performance using $V_{DD}$ scaling, a 40% frequency boost can be achieved with only $\sim110\%$ increase in leakage current, whereas a $\sim25\%$ performance increase using body-biasing results in 200% increase in leakage. The converse is true when reducing frequency, giving advantage to $V_{DD}$ scaling.

Nevertheless, leakage power can often be a negligible fraction of total power (in the case of the adder operating at $V_{DD} = 1.0V$, simulations show it is 1% of total power) and adjusting $V_{DD}$ impacts active power in a quadratic fashion ($P \propto CV_{DD}^2$). Thus, the impact of increased leakage currents is dependent on the ratio of leakage power to the total power. Eq. 4.7 gives the total power of a system, where the first term constitutes the active power of the circuit, the second term is the contribution of idle leakage currents, and $\alpha$ denotes the effective duty-cycle or percentage of time

(activity factor) that the circuit is actively operated.

$$P_{TOT} = \alpha C_{eff} V_{DD}^2 f + I_{leak} V_{DD} \qquad (4.7)$$

As a result, the required application throughput versus percent of time spent in idle modes will determine which of the two techniques, or even some combination of the two, is best used to mitigate variation. In high-performance domains where the activity factors are close to 1 and leakage power is insignificant relative to total power, body-biasing may be used effectively until the leakage power becomes a significant fraction of the total power, and body-biasing is advantageous as it does not, to first order, impact active power. Idle-mode dominated systems will likely benefit from increasing $V_T$ (by either changing the nominal $V_T$ offered by the technology or by appropriate body-biasing) and decreasing $V_{DD}$ to reduce idle power. Power-supply voltage scaling can be utilized to mitigate variation, as the increase in active power is not significant relative to the total power, and the increase in leakage energy is offset by the increase in nominal $V_T$ which decreases leakage energy.

Both techniques have associated overheads in terms of area, routing and design complexity when used in a within-die context, and typically have not been used in such a context. However, the use of body-biasing as a die-to-die variation mitigation technique, to lower idle power of "fast" die without severely impacting performance by biasing the entire substrate of a normal twin-well product, has been studied and commercially adopted [78, 120].

## 4.5 Summary

The beginning of this chapter advocated abstraction as a useful tool not only in digital circuit design but also in characterization and modeling of digital circuit variation. Specifically, variation models that can accurately model variation in digital gates or circuits and "plug" into static timing analysis tools frequently used by digital designers are of significant value. Many of these timing tools have also been advanced

to include spatial correlation within the analysis, but with little manufacturing data to validate this need. Moreover, mitigation techniques can be highly dependent on the amount of spatial correlation present.

Motivated by these needs, a test-chip was designed and implemented to characterize variation in common digital circuits, including the explicit ability to decompose the various spatial components of observed variation. Variation measurements highlighted the fact that averaging of random variation mitigates the impact of variation on larger gates and circuits. However, at lower voltages this becomes less true: the magnitude of variation increases as does the relative contribution from within-die components. Further analysis revealed that spatial correlation can often be confounded with systematic variation or even correlation arising from circuit structure. When these confounding components are properly accounted for, there appears to be no spatial correlation present in within-die variation. These results suggest that incorporation of spatial correlation in statistical static timing models is not necessary. Rather, incorporating positional dependencies (systematic components) may be more valuable, although quantifying these dependencies is likely difficult due to the apparent randomness of the systematic components.

Body-biasing and power-supply tuning were also evaluated as variation mitigation techniques. While both techniques are capable of mitigating even large amounts of variation, each requires different overhead in terms of die area, circuitry and impact to both active and idle power. Tuning power-supply to achieve variation mitigation has less impact on idle leakage currents than equivalent tuning using body-biasing but can also impact active power (both positively and negatively). Depending on the relative contribution of idle versus active power to the total power, either technique or a combination of these techniques might be appropriate. As a result of this analysis, the following chapter evaluates the efficacy of power-supply tuning to reduce the energy required to mitigate variation at an even higher level of abstraction, in multicore processor architectures.

# Chapter 5

# Variation Mitigation at the Architectural Level

Up to this point, our primary effort has been in characterizing, decomposing and quantifying the impact of variation on circuit performance. Understanding how to mitigate or cope with performance uncertainty is a necessary component in the formation of an integrated view of the variation picture. At the end of the last chapter, relatively simple mitigation schemes at the circuit level were evaluated; however, technology and process scaling have resulted in increasingly integrated modern designs which are effectively systems-on-a-chip. Large designs now incorporate multiple functional blocks and pack hundreds of millions to billions of transistors on a single chip. At these scales, it is no longer sufficient to focus on process variation from the device or even circuit perspective. Instead, the entire system is the appropriate context, as analyzing and optimizing individual circuits or blocks may result in sub-optimal solutions in the context of the system.

Narrowing down the scope of digital systems to analyze, this chapter focuses on massively multicore microprocessors, or Chip Multi-Processors (CMPs), which have hundreds or thousands of small, homogenous processing cores. Section 5.1 begins by analyzing the potential impact that variation can pose to these systems from a joint performance, energy and yield perspective. A mitigation scheme involving the addition of power-supply voltages is proposed in Section 5.2, and is followed by an ex-

planation of the simulation methodology utilized in evaluating the proposed scheme. The results of minimizing energy by adding power-supply voltages are presented and discussed in Section 5.3, along with details that must be considered when implementing such a scheme.

## 5.1 Impact of Variation on Multi-Core Processors

The trade-offs between mitigating performance variability and other key product metrics began to emerge in Section 4.4 of the last chapter: mitigating the impact of process variation by tuning either body-bias or power-supply voltage impacted idle and/or active power. In general, most proposed and implemented variation mitigation techniques involve trade-offs between die area, design complexity, power, cost and yield, making evaluation of any technique a difficult, multi-dimensional problem. The magnitude of these relationships is hard to appreciate at the level of an individual circuit due to its relatively small size and power consumption, and in the context of multi-core processors mitigating variation becomes non-trivial in scope and complexity.

Effectively managing variation at these levels is critical, as high-performance multi-core processors, in which power and variation are intricately linked, are expected to scale to many tens if not hundreds or thousands of cores per die. In such systems, core-to-core frequency variations will arise due to underlying process variation. Discussions with computer architects reveal that both architects and operating system designers value operating frequency homogeneity at the system level unless the cost of ensuring it is too high. If this is the case, solutions such as Factored Operating Systems, where each core runs an OS servlet and shields software from underlying core differences [121], and self-aware software capable of detecting a core's power and performance state through a variety of hardware sensors, have been proposed. However, the value of homogeneity is believed to be greater than software solutions, if core-to-core performance variation lies in the range of 20-50%.

To ensure homogenous core frequencies, a number of techniques might be em-

ployed, but all have trade-offs that must be considered:

1. *Increase device sizes.* While increasing device areas can reduce susceptibility to intrinsic variation sources such as RDF and LER, systematic variation sources (e.g., mask errors, CMP) can still result in variation. Furthermore, increased power is necessary to charge and discharge the larger capacitances, as they scale with device area. Larger total die areas also result in fewer die per wafer, increasing costs.

2. *Detect and correct.* This is a catch-all term for circuit or architectural techniques (e.g., Error-Correcting Codes used in SRAMs, Razor flip-flops [114]) that are able to detect logic or timing errors and either correct them or restart computation with adjusted conditions guaranteed to ensure correct completion. Many of these techniques are heavily used, requiring area and power overhead and increased design complexity, but potentially allowing for homogenous core frequencies.

3. *Asynchronous architectures.* Such architectures have long been proposed (e.g., [122]) but never commercially implemented. Additionally, they are still susceptible to variation, making guarantees of identical throughput between asynchronous cores difficult.

4. *Lower the clock frequency.* Slowing the clock frequency of each core to that of the worst-performing core guarantees homogoneity and is attractive due to its simplicity, but unnecessarily leaves considerable performance on the table.

5. *Increase the voltage.* Increasing the voltage of the entire system to that required by the worst-performing core, as in Figure 5-1(b), ensures maximum performance but incurs substantial power penalties.

6. *Individual core voltages.* Providing unique power-supply voltages, as pictured in Figure 5-1(d), can be close to an optimal solution from a power/energy perspective, but requires significant design and area overhead for inclusion of the necessary voltage regulators or DC-DC converters.

(a) CMP block diagram  (b) $N = 1$, $E = 1.23$  (c) $N = 2$, $E = 1.083$  (d) $N = N_{core}$, $E = 1$  (e) Voltage scale

Figure 5-1: Block diagram of a CMP with each core able to select from $N$ voltages and example core voltages, with normalized energies, to meet a performance constraint.

The above list is by no means exhaustive, but it is evident that any "solution" includes significant undesirable components, most often in the form of increased power dissipation. Bowman et al. showed that 31-53% additional power/energy is necessary to overcome the impact of process variation at the $50nm$ technology node, if voltage scaling is used to maintain performance over the nominal case of no variation [123]. In justifying the push to thousand-core processors, Borkar acknowledges that fine-grain power management is necessary to fit these processors within the desired power envelopes [124]. For design and power delivery simplicity, Borkar suggests using two voltage supplies such that a core operates at either a frequency, $f$, or $f/2$ and uses the lower voltage when operating at $f/2$.

Given the near certainty with which we can expect variation to impact a design of this magnitude, we will define the "variation-induced energy overhead" as the energy required over and above that when an ideal mitigation solution, such as individual core voltages, is used ($E_{ideal}$). Mathematically, this can be defined as $\frac{E - E_{ideal}}{E_{ideal}}$. Simulations on a RAW processor core [125] ported to the $45nm$ technology node show this can be 20% or more, depending on the amount of variation, if a single system voltage is simply scaled upward to account for the worst performing core. Though not as pessimistic as predicted by Bowman, this magnitude of power/energy overhead is large enough to warrant more efficient solutions.

# 5.2 Mitigation Strategy: Multiple Power-Supply Voltages

We tackle the combined power and variability problem in generic multi-core processors with the introduction of one or more additional power-supply voltages to the system (Figure 5-1(c)). Specifically, we go beyond Borkar's suggestion of two system voltages, as we focus on efficient selection of the optimal value of a vector of power-supply voltages whereby *each core* of a chip multi-processor (CMP) is assigned a single voltage from within the vector in order to minimize total chip energy while meeting performance (frequency) *and* yield constraints. In this section, we formulate an analytic approach and provide an efficient iterative algorithm to find good power-supply vector values. When the vector is composed of only two voltages ($N = 2$), we prove there is only a single optimum, and we provide an efficient mathematical formulation for finding the optimum. When the vector contains more than two voltages ($N > 2$), the algorithm utilizes the $N = 2$ case to find local optima. Quantitative bounds on the performance are also formulated, and we qualitatively show that our algorithm behaves according to the derived bounds.

## 5.2.1 Problem Formulation

Energy in a multicore processor is computed as the sum of the individual core energies plus any shared resources, as shown in Eq. 5.1, where $N_{Core}$ is the total number of cores on chip. The individual energy/operation of each core is shown in Eq. 5.2, where $C_i$ is the effective switched capacitance in the core and $T$ is the cycle time required to complete an operation.

$$E = \sum_{i=1}^{N_{Core}} E_i + E_{shared} \tag{5.1}$$

$$E_i = E_{dyn_i} + E_{leak_i} = C_i V_{DD}^2 + I_{leak_i} V_{DD} T \tag{5.2}$$

The minimum cycle-time achievable by a core is a function of power-supply voltage and can be expressed as in Eq. 5.3, where $K$ and $V_T$ are parameters determined by the

113

critical path in the design and are subject to variation. $\alpha$ is a technology-dependent parameter.

$$T = \frac{KV_{DD}}{(V_{DD} - V_T)^\alpha} \qquad (5.3)$$

Before continuing with further formulation of the problem, we list some assumptions. Specifically, we assume a high-performance CMP, leading to the following assumptions:

1. Die area and leakage energy are dominated by SRAM-based caches [76]. Since SRAM caches are typically on separate power-supplies and there exists a considerable body of work on leakage energy reduction in SRAMs, we exclude leakage energy, $E_{leak_i}$, from the analysis for the time being — it will be revisted in Section 5.3.4.

2. $E_{shared}$ is the energy of shared caches, I/O and other peripheral circuit blocks surrounding the processor cores. Many, if not all, of these blocks have their own power-supplies separate from the processor cores. As a result, modifications in how the cores are powered do not, to first order, affect this component and will be omitted in the following analysis.

3. $C_i$ has small enough variance over all cores that we can treat it as a constant. In general, on applications capable of massively parallel computation, the activity factor is likely similar over all cores. Furthermore, the total capacitance in a core is the sum of millions or more individual capacitances. When accounting for aggregate variation, summation typically reduces variation in $C_i$ due to averaging of random variation.

Our goal is to *minimize E* subject to both yield and performance constraints. In particular, we wish to minimize $E$ such that some percentage of the cores, $0 \leq y_o \leq 1$, in a CMP achieve a certain minimum frequency (maximum delay) of operation, $f_{min}$. Mathematically, this can be stated using the following definitions. A core is labeled "acceptable" if its frequency of operation for a given voltage is greater than or equal

114

to the constraint. Otherwise, it is unacceptable:

$$A_i \stackrel{def}{=} \begin{cases} 1, & \text{if } f_i(V_{DD}) \geq f_{min}; \\ 0, & \text{otherwise.} \end{cases} \tag{5.4}$$

The summation of $A_i$ over all cores gives the number of acceptable cores:

$$N_{Acc} = \sum_{i=1}^{N_{Core}} A_i \tag{5.5}$$

And finally, the number of acceptable cores divided by the total number of cores must be greater than or equal to the yield constraint:

$$\frac{N_{Acc}}{N_{Core}} \geq y_o \tag{5.6}$$

These constraints can be achieved by allowing each core to select its own minimum power-supply voltage, denoted by $V_{i,min}$, so that $f_i(V_{i,min}) = f_{min}$, as shown in Figure 5-1(d).[1] This is the case where we have as many supply voltages as we have cores, $N = N_{Core}$. However, as discussed above, this solution introduces substantial overhead. Instead, we can use a smaller number of power-supply voltages $(N \ll N_{Core})$, as shown in Figure 5-1(c), and attempt to minimize the energy in such a case.

As defined above, $V_{i,min}$ is the minimum voltage required for each core such that all cores operate at the desired frequency and provide a homogeneous view at the system level. Core-to-core frequency variation will result in a distribution of $V_{i,min}$'s, represented by $f(V_{min})$ as illustrated in Figure 5-2, which describes the probability that a core requires $V_{min}$ to operate at the desired frequency. The Cumulative Distribution Function (CDF), $F(V_{min})$, is the percentage of cores that can successfully use $V_{min}$. As such, using only $N$ voltages and combined with the assumptions mentioned

---

[1] While $V_T$ can be adjusted using body-biasing, this is generally less efficient than modifying $V_{DD}$ as seen in the last chapter (Section 4.4.3) $K$ is a technology-dependent parameter that is unmodifiable.

Figure 5-2: Example $V_{min}$ distribution and discretization for 1K-core RAW processor based CMP

above, the total energy we wish to minimize is:

$$E_N = CN_{Core} \left( V_1^2 F(V_1) + \sum_{i=2}^{N} V_i^2 \left[ F(V_i) - F(V_{i-1}) \right] \right) \qquad (5.7)$$

The equation above amounts to a discretization of the second moment of the distribution, because as $N \rightarrow \infty$ it reduces to $E_{ideal} = CN_{core} \int V^2 f(V) dV$. As we do not have an infinite number of voltages, we depict this in Figure 5-2 with the solid black curve representing $E_{ideal}$ for $N = N_{core}$, the cumulative energy[2] that would be required if every core were assigned its associated $V_{i,min}$. By using fewer than $N_{Core}$ voltages, the energy curve is discretized, with voltage placements (computed using the algorithm presented in the following section) at the positions indicated by the dashed lines. This is similar to approximating a continuous integral by partitioning the interval and using Reimann sums of finite subintervals (distance between dashed lines of the same color). However, in this case the $E_{ideal}$ curve is a lower bound, as each core must be provided a voltage greater than or equal to its minimum required voltage. Performing this discretization, we see that using only one voltage (red upper-pointing triangle) results in the greatest overhead (nearly 17% in this example), or worst approximation of the $E_{ideal}$ curve. Increasing the number of voltages results

---

[2]$C$ and $N_{core}$ are normalized out due to normalizing all energies to the case of $E_{ideal}$.

in successively better approximations, and with ten voltages (light blue squares), the energy overhead is tiny.

Using both Eq. 5.7 and Figure 5-2, we see that to meet a yield constraint, $y = y_o$, we should pick $V_N$ such that it satisfies $V_N = F^{-1}(y_o)$, where $F^{-1}(y)$ is the inverse cumulative distribution function. In this work, we use a Gaussian distribution function and its associated CDF due to relative ease of analysis compared to other distribution functions. Furthermore, despite $N_{core} < \infty$, we use a continuous rather than discrete distribution, as the ensuing math is made more tractable. Both approximations introduce only small error, as will be shown in Section 5.3.3.

## 5.2.2    Energy Minimization

Although Figure 5-2 included voltage placements that minimized total energy for the given number of voltages being used, we did not discuss how a particular vector of voltages is chosen. To choose a vector of voltages that minimizes energy, we use a very simple, but highly efficient iterative algorithm. We first present the algorithm and then qualitatively discuss the performance of the algorithm.

### Minimum-Energy Voltage Selection Algorithm

When choosing a vector, $V^*$, of $N$ voltages ($V_1 < V_2 < ... < V_{N-1} < V_N$), $V_N$ is chosen to meet a yield constraint as discussed above, so we need only choose the other $N - 1$ voltages. We propose the Minimum-Energy Voltage Selection (MEVS) algorithm shown in Algorithm 1. The algorithm begins with all $N$ voltages spaced uniformly, and iteratively solves for the optimal $V_i$ given that all other $V_j, j \neq i$ are equal to their previous values, until none of the $V_i$'s change by more than $\epsilon$ from one iteration to the next. Solving for a single $V_i$ while keeping all others constant is equivalent to solving the simpler case of two voltages ($N = 2$), detailed next.

With only two voltages in the system, Eq. 5.7 reduces to:

$$E_2 = C N_{Core} \left( V_1^2 F(V_1) + V_2^2 \left[ F(V_2) - F(V_1) \right] \right) \tag{5.8}$$

```
V* = distribute Vᵢ's uniformly;
Initialize V*old to 0;
while (Vᵢ − Vᵢ,old > ε) ∀i do
    V*old = V*;
    foreach Voltage Vᵢ do
        /* Solve for local optimal Vᵢ holding all Vⱼ≠ᵢ constant    */
        Vᵢ = FindOptimal(V*, i);
    end
end
```

**Algorithm 1**: MEVS algorithm


$V_2$ is picked apriori to meet the yield constraint as discussed in the previous section, so the problem is reduced to optimally choosing $V_1$, which may only take values $0 \leq V_1 < V_2$. There is only one optimal choice for $V_1$ and the proof of this is shown in Appendix A.1. This single optimum is also seen in Figure 5-3(a) where the total energy is plotted versus $V_1$. When $V_1$ is smaller than $min(V_{i,min})$ all cores must use $V_2$ to meet the performance constraint and hence the energy is maximum. However, as $V_1$ increases, more cores utilize $V_1$ rather than $V_2$ and the energy decreases. As $V_1$ continues increasing, the cumulative energy of the cores utilizing $V_1$ grows faster than the decrease in energy resulting from switching from $V_2$ to $V_1$, resulting in a minimum energy point.



(a) Two system voltages      (b) Three system voltages

Figure 5-3: Single optimum for both $N = 2$ and $N > 2$


Since there is no closed-form for the Normal CDF, it is written in terms of the

error function, expressed as a Maclaurin series:

$$F(x) = \frac{1}{2}\left[1 + erf\left(\frac{x-\mu}{\sigma\sqrt{2}}\right)\right] \tag{5.9}$$

$$erf(z) = \frac{2}{\sqrt{\pi}}\sum_{n=0}^{\infty}\frac{(-1)^n z^{2n+1}}{n!(2n+1)} \tag{5.10}$$

$$F(x) = \frac{1}{2}\left[1 + \frac{2}{\sqrt{\pi}}\sum_{n=0}^{\infty}\frac{(-1)^n\left(\frac{x-\mu}{\sigma\sqrt{2}}\right)^{2n+1}}{n!(2n+1)}\right] \tag{5.11}$$

We substitute Eq. 5.11 into Eq. 5.8 using a finite number of terms from the Maclaurin series (in practice, three or four terms achieves good accuracy). The derivative of the resulting polynomial is taken, resulting in another polynomial. The roots of the latter are determined using the *roots* function in MATLAB. All but one root are invalid, lying outside the range of interest or being complex.

In the multiple voltage case, Eq. 5.7 is used but only a single voltage is solved for at a time, and the derivative with respect to that voltage is used; all other terms are constant. While Eq. 5.7 is non-convex in general, we have empirically observed that there is only a single global optimum, as seen in Figure 5-3(b) for $N = 3$. We have also observed that though the proposed algorithm finds local optima, they are exactly (or very close to) the global optimum (see Figure 5-8).

**Algorithm Performance**

Since we are unable to prove global optimality in the general case, we attempt to bound the energy overhead. This overhead is the difference between the "ideal" energy when each core has its own voltage (Eq. 5.12) and the energy when we have discretized Eq. 5.12 using $N < N_{core}$ voltages (Eq. 5.13).

$$E_{ideal} = \int_{-\infty}^{\infty} V^2 f(V)dV \tag{5.12}$$

$$E_{discrete} = \sum_{i=1}^{N} \int_{V_{i-1}}^{V_i} V_i^2 f(V) dV \qquad (5.13)$$

The bounds on the overhead, and thus algorithm performance, are determined by how the $N$ voltages are distributed. In the simple case of uniform intervals between the voltages, the overhead can be bounded by $O(\frac{\sigma}{N}\sqrt{\mu^2 + \sigma^2})$ (not shown here due to space constraints). However, with more intelligent spacings, the bounds can be tightened to $O(\frac{\mu\sigma}{N})$ as shown in Appendix A.2. Intuitively, the spacings chosen are such that the size of each interval balances out the voltage cost ($V^2$ in Eq. 5.12) across all intervals. One could also attempt to balance the entire energy over all of the intervals (i.e., balance $V^2 f(V) dV$). However, since the ratio $\frac{f(\mu)}{f(\mu+\sigma)}$ is constant (or the area underneath $f(V)$ remains constant despite changes in $\sigma$), this will only change the value of the constants in the bound.

As seen in Figure 5-4, the actual performance of the algorithm closely fits the derived bounds. In the case of Figure 5-4(c), the fit is quadratic rather than linear as was derived. Conceptually, this makes sense as we did not handle mean shifts in the derivation in Appendix A.2. If the mean shifts, both $E_{ideal}$ and $E_{discrete}$ increase in a quadratic fashion due to the $V^2$ term, and so the bound should also include a $\mu^2$ term as the actual performance indicates.

## 5.3 Mitigation Results

The power of the MEVS algorithm is in both the analytic formulation and the computational efficiency it offers. These characteristics allow for a fully analytic characterization of the energy overhead and the potential energy reduction based only on variation statistics ($\mu$ and $\sigma$). This section first focuses on such analysis to explore the possibilities of using additional voltages to both mitigate variation and reduce energy. It then details a custom simulation methodology employed to demonstrate energy reduction in a one-thousand core CMP. Since a one-thousand core CMP does not yet exist, we ported a RAW microprocessor core design to a $45nm$ technology and characterized the variation that might be observed in such a system.

(a) Overhead vs. N



(b) Overhead vs. $\sigma$



(c) Overhead vs. $\mu$

Figure 5-4: Actual algorithm performance (stars = simulation data points, dashed line = fit line)

## 5.3.1 Analytic Energy Reduction

With an analytic framework in place, the "variability-induced energy overhead" can be characterized. In Figure 5-5(a), the energy overhead for $+3\sigma$ yield is plotted versus the magnitude of variation present in the $V_{min}$ distribution $(\frac{\sigma}{\mu})$ and the number of voltages in the system $(N)$.[3] The linear dependence on $\sigma$ is again noticed, but focusing on $N = 1$, it is apparent that even for a modest amount of variation ($\sim 5\%$) the energy overhead is significant, approaching 30%. Given the measured results from the last chapter and the rougly linear dependence between delay and power-supply voltage observed in those measurements (see Figure 4-21(a)), at the $90nm$ node within-die

---

[3]For this analysis $\mu = 1.0V$ as this is the nominal voltage found in most state-of-the-art systems.

variation results in $\frac{\sigma}{\mu} \approx 1\%$ and consequently a 5% energy overhead. While fairly manageable at the $90nm$ node, the magnitude of variation typically increases and as the number of cores increases, the yield constraint will increase to perhaps $4\sigma$ or $5\sigma$, resulting in much larger energy overheads; the effects of changing the yield constraints are discussed and quantified below.



(a) $E_{overhead}$        (b) Reduction in $E_{overhead}$

Figure 5-5: Analytic computation of $E_{overhead}$ and reduction of $E_{overhead}$ when using additional voltages selected by the MEVS algorithm.

Looking at Figure 5-5(b), we notice that the amount of reduction in energy overhead is constant with the magnitude of variation. However, since the magnitude of the overhead is linearly increasing with the magnitude of variation, the energy reduction relative to the total energy will also increase linearly with variation. Furthermore, the addition of only a single new power-supply voltage ($N = 2$) provides the largest incremental energy savings no matter the magnitude of variation, with asymptotic energy reduction afterward.

## 5.3.2 Simulation Methodology

To test the MEVS algorithm and demonstrate the energy savings of using multiple system voltages in a real design, we used the RAW core, developed at MIT, as it was specifically developed for multicore applications [125]. However, the 64-core RAW processor was implemented in a mature $0.18\mu m$ technology node, requiring that the core be ported to a more leading-edge process before simulations could be carried

The diagram contains the following labeled elements:

**Within-Die Process Variation Model** — Probability vs Process Parameter (L, $V_T$)

**Die-to-Die Process Variation Model** — Probability vs Process Parameter (L, $V_T$)

**1K-point Monte Carlo Hspice voltage sweeps**

**Delay vs Voltage Curves** — Delay vs $V_{DD}$

$$D = \frac{KV_{DD}}{(V_{DD} - V_T)^\alpha}$$

Fit delay curves using MMSE

**Joint $V_T$, K Distributions** — Probability vs $V_T$ and K

**MATLAB 10K-point Monte Carlo analysis using analytic delay equation**

Probability vs Delay

Probability vs $V_{min}$

**MEVS Algorithm to determine optimal voltages**

Figure 5-6: Simulation methodology to efficiently evaluate MEVS algorithm and resultant energy savings.

out. This involved re-synthesis in Synopsys Design Compiler with a non-optimized predictive $45nm$ technology (PTM [113]) using FreePDK45, in combination with Nangate's OpenCell standard cell library [126]. As neither the FreePDK45 nor Nangate's standard cell library included a memory compiler, SRAMs were not implemented or included in any of the subsequent simulations; however, this is consistent with all of the above analysis where SRAMs were also excluded. Although not optimal, processor register files were synthesized using available standard cell flip-flops.

Static timing analysis using Synopsys PrimeTime was then performed to choose 20 critical paths for detailed further analysis, as depicted in Figu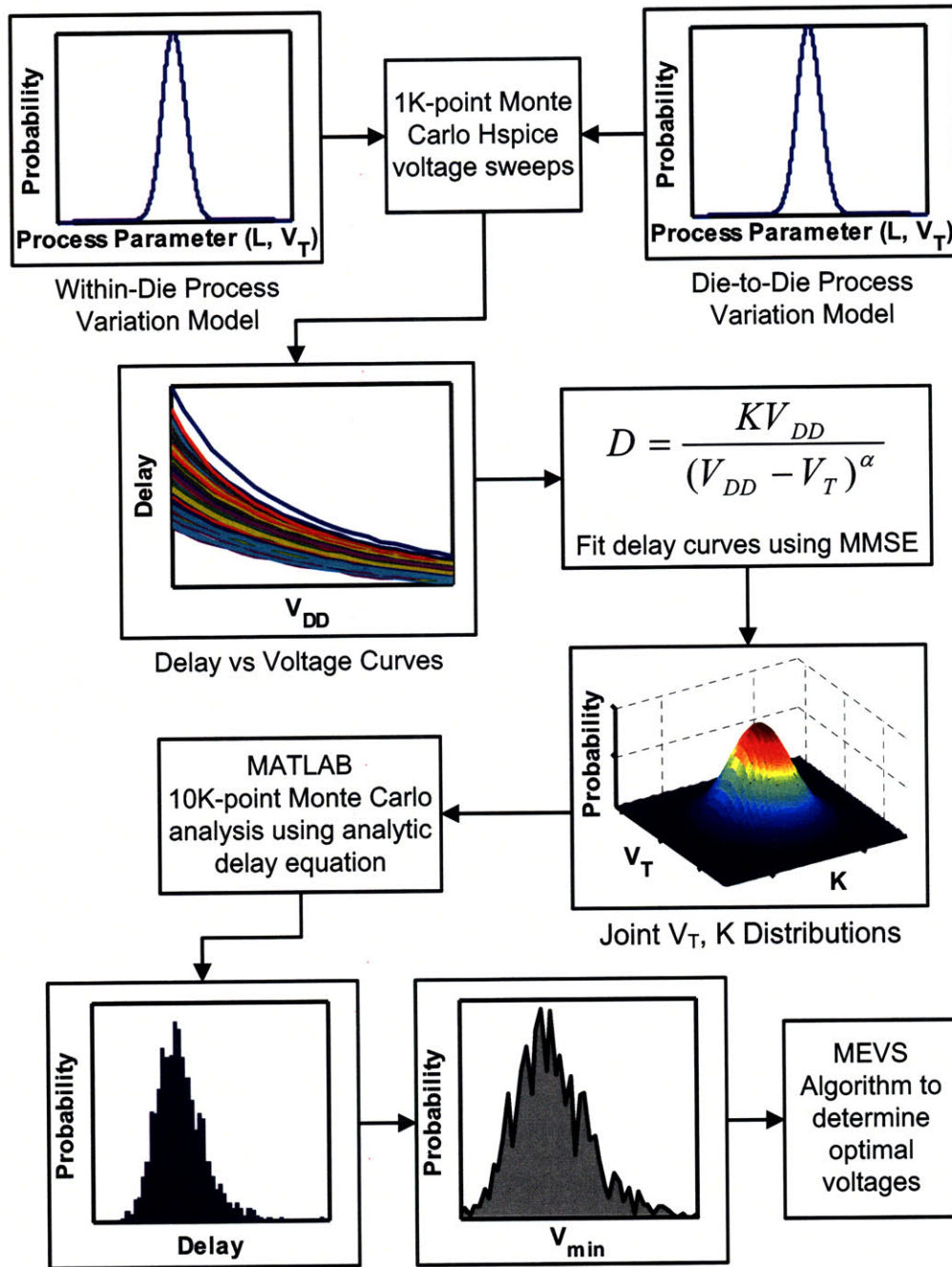re 5.3.2. A 1K-point Monte Carlo voltage-sweep analysis was done in HSPICE for each critical path, in order to analyze the effects of within-die random variation on each path. Sweeps for a single critical path are shown in Figure 5-7(a) with the associated probability distribution in Figure 5-7(b), showing increased delay variability as $V_{DD}$ decreases, consistent with the results from the last chapter. Within-die systematic variation was not modeled, as Section 4.3.2 of the last chapter indicated that random variation dominates in regular, arrayed structures.

Each of the 1K voltage-sweeps for each path was fit to the delay model in Eq. 5.3, and variation in the delay model parameters ($K$ and $V_T$) was characterized. In particular, for each voltage sweep of each path, the appropriate values of $K$ and $V_T$ were determined by achieving the best fit[4] to Eq. 5.3. Doing this for all 1000 sweeps resulted in a distribution of $K$ and $V_T$ for each path. The statistics of each distribution as well as any cross-correlation between $K$ and $V_T$ (generally low) were subsequently computed and saved for further simulation. For this process, the fit values of $K$ and $V_T$ had between $2 - 4\% \frac{\sigma}{\mu}$ variation for each path.

Once this characterization was complete, the statistics of the $K$ and $V_T$ distributions were used to model 10K one-thousand core CMPs using Eq. 5.3. Analytic modeling allowed efficient simulation of many more paths, ten million in this case, than would otherwise be achievable using time-domain simulation. For each core,

---

[4]The common Minimum Mean Square Error (MMSE) methodology is used for fitting non-linear data.

(a) Critical path delay versus $V_{DD}$

(b) Probability distribution of critical path delay versus $V_{DD}$

Figure 5-7: 1K-point Monte-Carlo voltage sweeps for a single critical path

random values for $K$ and $V_T$ were generated and used to generate the delay versus $V_{DD}$ curves for each critical path. Each CMP was also subjected to a zero-mean, normally distributed mean ($\mu$) shift to simulate die-to-die variation. These curves were then combined to determine the minimum supply voltage, $V_{i,min}$, required to meet a user-defined delay constraint. This allowed generation of the $f(V_{min})$ distribution for each CMP to which the MEVS algorithm was applied.

### 5.3.3 Results & Analysis

The above simulation methodology, applied on 10K samples of a one-thousand core multicore processor based on the RAW core, allowed efficient evaluation of the energy savings as a result of adding voltages to the system.

**Effect of Approximations**

To understand how the approximations mentioned above affect optimal voltage selection and the resultant energy reduction, we compared a subset of the 10K voltage vectors selected by the MEVS algorithm to the actual optimal vector in each case. Since finding the actual optimal vector requires exhaustive search over a large multi-dimensional space, time and computational constraints limited this comparison to no

more than six voltages.

Firstly, Figure 5-2 shows an example $V_{min}$ distribution in a single CMP. Despite not matching any common distribution, the bulk (80%) of the distribution can be approximated with the Gaussian distribution, resulting in only small error. More importantly, Figure 5-8 shows that the MEVS-selected vector is very close to the globally optimal vector despite using: 1) a Gaussian distribution as an approximation to the actual distribution, 2) the truncated Maclaurin series approximation for the Gaussian CDF, 3) continuous rather than discrete math, and 4) slightly larger $\epsilon$ as $N$ increases to aid in convergence time.[5]



Figure 5-8: Vector distance and energy difference between MEVS-selected and globally optimal vectors.

Even with six voltages in the system, the average vector distance from optimal is $10mV$, within the ripple of any power-supply voltage (typically no less than $10mV$). As the number of voltages increase, an increase in the distance from optimal is expected simply due to adding the error of each voltage. However, even in the worst-case of a $25mV$ distance, this implies that the five additional voltages are on average no more than $5mV$ away from their optimal values. More importantly, these small distances from the optimal values result in energy differences between the MEVS-selected vectors and the globally optimal vector of much less than 1%.

---

[5]Using many voltages in a system would likely not be practical, and for $N \leq 5$ keeping $\epsilon$ small ($\sim 1mV$) has no effect on convergence time.

## Energy Reduction

With the effect of approximations shown to be small, we next analyze energy reduction in the RAW core as a result of using multiple voltages. Figure 5-9 plots both the total energy reduction (energy difference between using a single voltage and multiple voltages), as well as the reduction in the amount of variability-induced energy overhead (energy over and above each core having its own voltage $V_{i,min}$). By adding just a single additional voltage ($N = 2$), anywhere between 59-75% of the variability-induced energy overhead is eliminated, resulting in a total energy savings of 6-16%, with an average savings of 9%. These results are in-line with the analytic results above, and are expected from observed variability in the delay and $V_{min}$ distributions of $\sim 3\%$. The large range of energy savings noticed is due primarily to die-to-die variation, which results in mean shifts of the $V_{min}$ distribution as opposed to greater magnitudes of within-die variation.



Figure 5-9: Energy reduction

The addition of more voltages does increase the energy savings but at diminishing returns: with five voltages, roughly 90% of the energy overhead is eliminated, and it would take 995 additional voltages to reach $E_{ideal}$. Since the absolute energy savings are dependent on the magnitude of variation, as CMPs are scaled to smaller processes where variation is expected to increase, the energy overhead of using only a single voltage for all cores will also increase, as seen in Figure 5-5(a). Use of our

voltage selection algorithm will result in larger absolute energy savings as the relative magnitude of variation $\left(\frac{\sigma}{\mu}\right)$ increases.

**Effect of Yield-Constraint Choice**

In the above analysis, the yield constraint was such that the last voltage in the system, $V_N$, had to accommodate the worst-performing core, or stated differently, all cores had to function at the required frequency. However, in such a massively parallel system, having a 100% yield may not always be necessary, nor may it be efficient from a performance-energy perspective. To quantitatively analyze this, an appropriate metric is required. Since performance and energy are both individually important, a metric that includes both is used: for this analysis, our metric is the ratio of the total performance of the system to the total energy in the system.

The total performance in the system is proportional to the product of clock frequency ($f$), number of instructions completed per clock ($IPC$), and number of operational (or yielding) cores ($y_o \times N_{core}$). The total energy in the system is given by Eq. 5.7 multiplied by the yield, $y_o$, and the frequency, $f$. Since two system voltages provide the most incremental benefit, this analysis is only performed for $N = 2$, and so Eq. 5.8 is used in place of Eq. 5.7 resulting in:

$$\frac{\text{Tot. Perf.}}{\text{Tot. Energy}} = \frac{y_o N_{core} IPC f}{y_o E_2 f} \tag{5.14}$$

$$\frac{\text{Tot. Perf.}}{\text{Tot. Energy}} = \frac{IPC}{C\left(V_1^2 F(V_1) + V_2^2 \left[F(V_2) - F(V_1)\right]\right)} \tag{5.15}$$

Although it would seem that Eq. 5.15 has no dependence on number of yielding cores, recall that there is an implicit dependence, as $V_2$ is selected a priori such that it meets the yield constraint ($V_2 = F^{-1}(y_o)$), and so there is still a dependence on yield. Furthermore, for the purposes of this analysis we assume that both $IPC$ and $C$ are constant (or do not change significantly per core), and are simply scaling factors that can be removed from the analysis.

When this analysis is performed, the combined performance/energy metric is

roughly constant, as seen in Figure 5-10(a). This is expected, as both performance and energy scale linearly with the number of operating cores ($y_o N_{core}$). However, there is a slight decreasing trend as the energy does not strictly scale linearly with yield, reflecting the necessary increase in both $V_1$ and $V_2$ to support additional poor-performing cores. As the yield constraint increases to roughly 85-90%, the increase in $V_2$ accelerates due to the exponential nature of the tails of the distribution as seen in Figure 5-10(b), resulting in faster decreases in the performance/energy metric as shown in the inset of Figure 5-10(a).



(a) Performance/Energy metric

(b) Required $V_1$, $V_2$

Figure 5-10: Joint performance/energy metric versus yield

Looking at the incremental change in performance relative to the incremental change in energy with increasing yield constraint, the right axis (green plot) of Figure 5-10(a) shows that the change is relatively constant until $y_o \geq 85\%$, where it sharply decreases. Intuitively, this means that for a constant increase in energy, a constant increase in performance is achieved until $y_o \geq 85\%$, at which point the incremental increase in energy required for the same incremental gain in performance becomes increasingly large. This result suggests that there may be an upper-bound on practical yield constraints unless total computational throughput is of the essence; a similar conclusion was reached in [95] where energy reduction is achieved solely by turning off power-hungry cores.

Another way to arrive at the same conclusion is to explore the energy overhead as a function of the desired core yield. Figure 5-11(a) shows an analytic computation

of $E_{overhead}$ as a function of the yield for a fixed amount of variation in the $V_{min}$ distribution[6] ($\frac{\sigma}{\mu} = 3\%$ in this case, according to the observed variation in the RAW core). The energy overhead increases exponentially with the desired yield constraint, and so reducing the yield constraint by a few percent can have a large impact on energy. Nevertheless, even if the yield constraint is reduced to 90%, $E_{overhead} \approx 8 - 10\%$, which is still significant enough to warrant adding voltages to the system. Figure 5-11(b) shows that a single additional voltage ($N = 2$) still provides the most incremental reduction in $E_{overhead}$; however, the benefit is somewhat reduced as the yield constraint is decreased.



(a) $E_{overhead}$

(b) Reduction in $E_{overhead}$

Figure 5-11: Analytic computation of $E_{overhead}$ and reduction of $E_{overhead}$ versus core yield constraint.

## 5.3.4 Practical Considerations

Modern microprocessors have many power/performance modes as well as other design constraints that must be considered when attempting to implement a multiple voltage system. The following are some of the more salient aspects of physical systems, with a brief analysis of how each affects the framework and results above.

---

[6]The number of voltages is intentionally limited to $N \leq 3$ for both plots in Figure 5-11, as the MEVS algorithm has difficulty assigning voltages due to the limited distance between minimum and maximum $V_{min}$, especially as the yield constraint is reduced. In this regime, the Maclaurin approximation is not sufficient to properly model the very narrow minimum (i.e., the derivatives change too rapidly in the vicinity of the minimum).

## Integration of Leakage Energy

The analytic framework above ignored leakage energy, because leakage energy typically has been small relative to the active energy of high-performance systems. However, there is evidence that this is changing, with gate and sub-threshold leakage currents growing [127], and so it may be advantageous to include leakage energy in the above formulation.

To do so, the leakage component in Eq. 5.2 must be included, where $I_{leak_i} = KV_{DD_i}e^{\eta V_{DD_i}}$. In general, $K$ and $\eta$ are process constants that vary from transistor to transistor. In this analysis, these parameters can be thought of as effective values for an entire core. Importantly, the averaging that occurs when many transistors are aggregated reduces the magnitude of variation in total core leakage current, as shown in Figure 5-12, where the relative variation in core leakage currents is inversely proportional to the square root of the number of transistors in the core. Consequently, $K$ and $\eta$ can be thought of as process-determined constants.
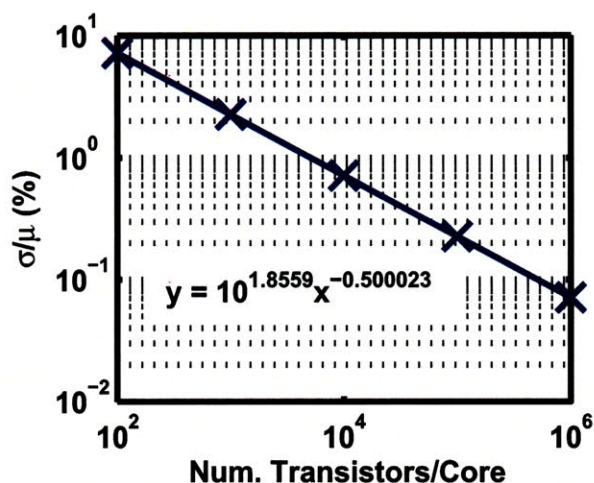


Figure 5-12: Variation in core leakage currents reduces to $\ll 1\%$ with even moderate core size.

In the two-voltage case, the addition of a leakage energy term results in the following reformulation:

$$E_2 = \left(CV_1^2 + KV_1e^{\eta V_1}\right) F(V_1) + \left(CV_2^2 + KV_2e^{\eta V_2}\right) [F(V_2) - F(V_1)] \qquad (5.16)$$

131

Taking the derivative of this reformulated energy to find the optimal $V_1$ results in the following constraint, analogous to Eq. A.1 found in Appendix A.1:

$$\frac{f(V_1)}{F(V_1)} = \frac{2V_1 + (1+\eta)\frac{K}{C}V_1 e^{\frac{\eta V_1}{nV_{th}}}}{(V_2^2 - V_1^2) + \frac{K}{C}\left(V_2 e^{\frac{\eta V_2}{nV_{th}}} - V_1 e^{\frac{\eta V_1}{nV_{th}}}\right)} \qquad (5.17)$$

These additional terms relative to Eq. A.1 do not change the fact that there remains only one optimal solution for $V_1$. The reasoning behind this is identical to that when leakage energy is not considered (Appendix A.1): despite the additional terms, the right-hand side of the equation still monotonically increases with $V_1$. However, these additional terms require knowledge of the effective switched capacitance ($C$), total number of gates and additional process-dependent terms, in order to compute $K$ and $\eta$. Some of these factors (particularly $C$) require extensive simulations on a highly optimized design to be meaningful. Owing to the highly *non-optimized* netlist generated by synthesis with a non-optimal library for our multicore case study, further analysis including leakage energy has not been performed, but should be a subject of future work.

## Temperature

The power dissipation associated with adjusting voltages also impacts the local temperature on the die. Temperature also impacts leakage energy significantly as the leakage current is exponentially dependent on threshold voltage ($V_T$) and the thermal voltage ($V_{th} = \frac{kT}{q}$), which are both affected by temperature. This can be seen in Figure 5-13, where the change in leakage current is plotted versus temperature and $V_{DD}$, and the exponential dependence on both parameters is visible.

Humenay et al. show that if voltage-scaling alone is used to compensate for the impacts of variation, increases in both leakage and temperature necessitate either more expensive cooling solutions or thermal throttling, leading to dynamic performance asymmetries [67]. Compared to this worst-case of scaling a single voltage system, adding voltages significantly reduces active power as demonstrated above,

Figure 5-13: Leakage current as a function of $V_{DD}$ and temperature (normalized to $V_{DD} = 1.0V$ and $T = 25°C$)

but also has a more marked impact on leakage power above than when considering DIBL alone: local temperatures will also decrease due to decreased active power dissipation resulting in lower leakage power as well. To quantify this effect, knowledge of the thermal characteristics of the cooling solution being utilized is necessary. In particular, the thermal impedances can be incorporated into leakage models which fit into the revised energy calculation above, and can be used to better select the optimal voltages based on temperature as well as duty cycles.

**Impact of Memory Subsystems**

Up to this point, we have implicitly ignored the impact of memory subsystems, as many techniques have been applied to reduce their active power, including putting large portions of the cache into sleep modes. For example, the Intel Dual Core Xeon processor features 18MB of L2 and L3 caches, but only 0.08% of the caches are actively powered for a given cache access. The large caches utilized on this processor allows the cache to be organized into many smaller arrays and sub-arrays, which can individually be put into sleep states. This results in 0.75W/MB average power for the caches, leading to the caches contributing less than 10% of the overall power budget [128].

However, in massively parallel multicore systems, each core will likely have much

smaller caches, in which a larger fraction of the cache will be actively powered for a cache access. More significantly, leakage of the many un-accessed lines in the active portion of the cache will result in an increase of the fraction of total power associated with the memory subsystems. In the case that there are additional globally-shared caches, these memory systems will increase the shared energy component relative to the total energy (see Eq. 5.1). Though this component must be incorporated, caches are typically operated on separate power-supplies from the computational cores and will likely not affect the optimization of the core power-supply voltages. The shared energy may affect how many and which cores are turned off; some preliminary work evaluating this has been done in [95], in which a core is turned off if the incremental power required to operate it is greater than the power of the shared blocks, amortized over all operating cores, if that core were not turned on.

## Frequency Scaling Systems

To conserve energy when application throughput requirements are not maximal, many digital systems employ frequency scaling, which also allows for scaling the system supply-voltage to further reduce both active and leakage energy. These systems are typically referred to as Dynamic Voltage-Frequency Scaling (DVFS) systems. If such systems are to also use multiple system voltages, an important question is whether each core retains its relative voltage assignment as the system voltages are scaled. For example, if a core is assigned to $V_i$ at a high-frequency/voltage setting, will it remain assigned to $V_i$ at a lower frequency/voltage setting? Or, is it necessary to recompute the core assignments for every frequency/voltage setting?

Intuitively, though the delay versus voltage curves observed in Figure 4-21 are roughly linear, the coefficients for each core (adder) vary and moreover, the computation of optimal voltage is a non-linear operation. As a result, the expectation is that there are some cores whose $V_{min}$ is relatively close to a given $V_i$ and will be reassigned to either $V_{i-1}$ or $V_{i+1}$ when optimal voltages are recomputed for a different frequency/voltage setting. This expectation is born out in simulations of a two voltage ($N = 2$) system shown in Figures 5-14 and 5-15, where the upper triangles
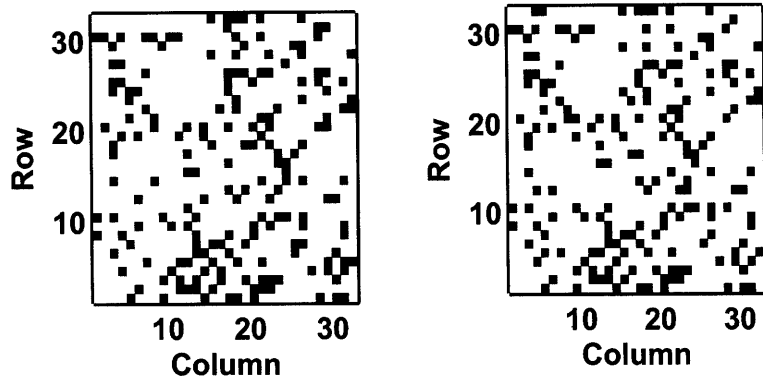
(red) indicate cores that have changed from using $V_1$ to $V_2$, and lower triangles (blue) indicate cores that have changed from $V_2$ to $V_1$, as the frequency and voltage change from a nominal frequency of $1.25GHz$ with $V_1 = 1.025V$ and $V_2 = 1.084V$ to a new frequency of $625MHz$.

These results indicate that, in addition to recomputation of optimal voltage, many cores must be reconfigured to use a different power-supply voltage. Practically, this can be implemented by pre-characterization at distinct frequency/voltage settings during sort/test and populating a ROM or other look-up table which contains this configuration information. By making use of the fact that frequency and voltage are roughly linearly related, this need not be done for every frequency/voltage setting. Rather, the end-points could be characterized and linear interpolation used (with some small margin to account for the slight non-linearities present).

Alternatively, it is also possible to use the same relative core-voltage configurations and pay an energy penalty. If $L_{F1}$ is defined as the original set of cores using $V_1$ at frequency $F_1$, and $H_{F1}$ is the remaining set of cores using $V_2$, then $V_{1_{Fi}}$ is set to $max(V_{min}(F_i))\forall L_{F1}$ and $V_{2_{Fi}}$ is picked as before. Note that the $V_{min}$ distribution is a function of the desired frequency setting (i.e., $V_{min}(F_i)$). When this approach is taken, the reduction in energy overhead is smaller relative to using the appropriate configuration, as summarized in Table 5.1. When the frequency and voltage are close to nominal (within 20%), the penalty of using the nominal core-voltage configuration is small. As the magnitude of frequency/voltage scaling increases, the penalty quickly grows: at $F_i = 0.5F_{nom}$ the penalty is a 15% loss in potential energy overhead reduction.

**Voltage Regulator & DC-DC Converter Efficiencies**

In the case where each core is provided with its own power-supply, a single system voltage is distributed over the entire die and linear voltage regulators are required to convert from the system voltage to the core-specific voltage, $V_{i,min}$. As mentioned previously, voltage regulators incur not only an area overhead, but an energy overhead as well due to an efficiency less than 100%: the efficiency of a voltage regulator is $\frac{V_{out}}{V_{in}}$

(a) $F = 1.250GHz$, $V_1 = 1.025V$, (b) $F = 1.125GHz$, $V_1 = 0.972V$, $V_2 = 1.084V$ $V_2 = 1.026V$

(c) $F = 1.000GHz$, $V_1 = 0.920V$, (d) $F = 0.875GHz$, $V_1 = 0.872V$, $V_2 = 0.970V$ $V_2 = 0.918V$

(e) $F = 0.750GHz$, $V_1 = 0.827V$, (f) $F = 0.625GHz$, $V_1 = 0.783V$, $V_2 = 0.871V$ $V_2 = 0.826V$

Figure 5-14: Core voltage assignment patterns. Dark squares indicate core is using $V_2$ while white squares indicate $V_1$.

(a) $F = 1.125GHz$, $V_1 = 0.972V$, (b) $F = 1.000GHz$, $V_1 = 0.920V$, $V_2 = 1.026V$ $V_2 = 0.970V$



(c) $F = 0.875GHz$, $V_1 = 0.872V$, (d) $F = 0.750GHz$, $V_1 = 0.827V$, $V_2 = 0.918V$ $V_2 = 0.871V$



(e) $F = 0.625GHz$, $V_1 = 0.783V$, $V_2 = 0.826V$

Figure 5-15: Pattern of cores changing between $V_1$ and $V_2$ as voltage/frequency settings are changed.

| Frequency | $E_{Overhead}$ Reduction | |
| :---: | :---: | :---: |
| | Ideal Core-Voltage Config. | Nominal Core-Voltage Config. |
| $1.250GHz$ (Nominal) | 65.53% | 65.53% |
| $1.125GHz$ | 65.56% | 64.79% |
| $1.000GHz$ | 65.58% | 62.81% |
| $0.875GHz$ | 65.56% | 59.59% |
| $0.750GHz$ | 65.48% | 55.33% |
| $0.625GHz$ | 65.30% | 50.11% |

Table 5.1: Energy penalty of adjusting $V_1$ to match nominal core-voltage configuration.

and $V_{in}$ is always larger than $V_{out}$. $V_{in}$ in this case will be the maximum core voltage required, $max(V_{i,min})$, and $V_{out}$ is simply the desired voltage, $V$. The ideal energy specified by Eq. 5.12 must now be modified:

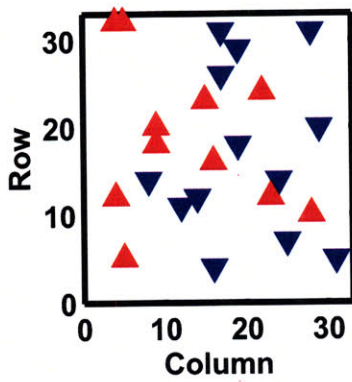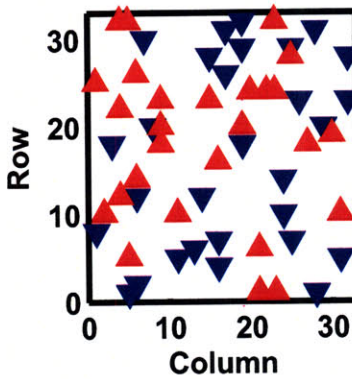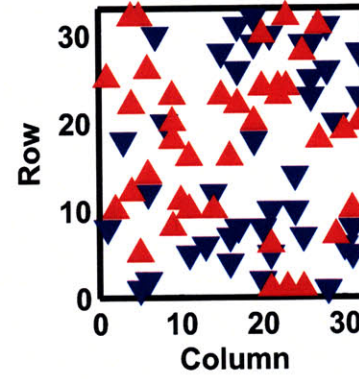$$E_{ideal} = \int_{-\infty}^{\infty} V^2 \frac{max(V_{i,min})}{V} f(V) dV = max(V_{i,min}) E[V] \qquad (5.18)$$

This modification means that $E_{ideal}$ is now greater than when these inefficiencies were ignored. In Figure 5-16, $E_{ideal}$ is computed according to the new definition in Eq. 5.18. The energy when using fewer voltages relative to this new definition is under 100% for $N > 1$, indicating energy *lower* than $E_{ideal}$ when accounting for voltage regulator inefficiency. Thus, having substantially fewer voltages than one for each core ($N << N_{core}$) results in improved energy reduction relative to each core having its own voltage, especially as the magnitude of variation increases and $max(V_{i,min})$ becomes larger and larger, decreasing efficiency further. In some sense, the ideal case, when accounting for these inefficiencies, is no longer ideal.

Of course, generation of the $N$ system voltages using DC-DC converters also incurs a small energy penalty due to sub-100% efficiency. However, DC-DC converters are typically off-chip[7] and have efficiency curves that are only dependent on load rather than any relationship between input and output voltages. Since a DC-DC converter is always used to generate the die voltage (i.e., even when $N = 1$), splitting a single

---

[7]Recently, there has also been analysis in integrating per-core high-frequency DC-DC converters on-chip showing energy savings when combined with DVFS [129]. Since the analysis was only performed for a small number of cores (4), it remains to be seen whether this is scalable to a much larger number of cores.

Figure 5-16: Energy using multiple system voltages relative to ideal energy when voltage regulator efficiencies are taken into account.

DC-DC converter into multiple converters does not, to first order, change the energy analysis.[8] There is, however, a board area/component trade-off as multiple converters, requiring typically large inductors and capacitors, are now required.

## Power Transistors & Routing Resources

Lastly, each voltage requires global routing resources for power distribution, but this can be mitigated in the case of two voltages to a large degree by reducing the width and density of each voltage's power grid since power will be divided roughly equally over the voltages. More importantly, power-multiplexing transistors are required. These transistors are necessarily large to handle the relatively large currents necessary to power each core. However, they can also be used to power-gate the entire core, as Intel has done in their most recent Nehalem architecture to eliminate leakage power when the core is unused [127].

---

[8]While the loads do change and DC-DC converter efficiencies are a function of load, the DC-DC converters can be optimized for the expected loads based on the number of system voltages implemented.

## 5.4 Summary

This chapter made clear that performance, power/energy, cost and variation are intricately linked to each other. As system architects value core homogeneity in multicore processors, the burden of finding *energy-efficient* variation mitigation solutions increases. Balancing other constraints, such as area overhead and design complexity, must also play into any mitigation technique. Evaluation of some of the more common mitigation schemes, in the context of these constraints, revealed a need for an energy-efficient technique capable of guaranteeing both performance and yield constraints, while introducing minimal overhead.

Introducing additional system voltages (fewer than one per core) provides a compromise between excess design/area overhead and energy efficiency while maximizing performance. An analytic framework capable of optimal voltage selection to minimize energy was developed and analyzed in-depth. Specifically, a simple, efficient algorithm to select optimal voltages formed the basis of this framework, allowing for multiple types of analysis. Though the optimization problem was not convex in general and simplifying approximations were used, the algorithm is nevertheless able to find optimal (for $N = 2$) or near optimal solutions. Furthermore, the behavior of the algorithm was mathematically bounded and shown to perform according to the bounds.

Using a custom simulation methodology and the MEVS algorithm, a core designed specifically for multicore contexts was simulated, to observe the magnitude of performance variation and the impact of introducing additional voltages to a massively parallel 1K-core processor. Analysis showed that a single additional power-supply voltage provides the greatest incremental impact, with 59-75% reduction in the "variation-induced energy overhead" and 6-16% total energy reduction. The desired yield constraint also has a significant impact on energy reduction: though counterintuitive, turning off a small fraction of the thousand cores can provide a positive trade-off between performance and energy, with multiple system voltages further improving this trade-off. Lastly, when voltage regulator efficiencies are properly ac-

counted for, using only a few system voltages ($2 \leq N \leq 10$) provides greater energy reduction than supplying each core with its own voltage.

# Chapter 6

# Thesis Summary & Future Work

This thesis has shown that process variation plays a significant role in area, power, performance, design complexity and cost of VLSI devices, circuits and systems. Beginning at the device level and moving upwards to circuits, we demonstrated that appropriate characterization and decomposition of process variation is critical to robust design as well as effective mitigation of the negative impacts of variation. This concluding chapter summarizes this thesis and presents ideas for follow-on work, since despite the many advances in understanding and mitigation of process variation, there is still much room for improvement.

## 6.1   Thesis Summary

With continued scaling of transistor dimensions, the magnitude of process variation has, for the most part, increased. Though the absolute variation has been kept under control (i.e., $\sigma$ reduces from one technology generation to the next), the relative magnitude ($\frac{\sigma}{\mu}$) has tended to grow. The sources of this variation are themselves varied, ranging from limited process module controllability as in sub-wavelength lithography, to fundamental sources of variation like random dopant fluctuation. Each variation source has its own unique signature, and thus characterization and modeling of these variation sources is most accurate when performed at the level of fundamental parameters. A key aspect determining the signature is the spatial decomposition of

the particular variation source. Since the spatial component of variation can influence both process development and product design, characterizing this component appropriately is increasingly critical.

Using sub-threshold current measurements to isolate threshold voltage variation from other sources of variation, we found that the magnitude of variation in $V_{T_o}$ is dependent only on the area of a device and, for a given device area a Gaussian distribution best describes the observed statistics. Due to further analysis showing a lack of both within-die spatial correlation and any systematic components, we showed that variation sensing and mitigation are highly dependent on the region of device operation. Low-power operation requires very different sensing techniques than high-performance operation, in particular the need to sense and measure *within* a critical circuit rather than *near* it. This potentially increases design complexity, as care must be taken to ensure that accurate within-circuit performance measurement does not significantly degrade performance or power, especially if the system must scale between operating modes.

This is evident again at the circuit level, where relative variation in circuit performance increases substantially as $V_{DD}$ is lowered. However, an all-digital, in-circuit technique to measure variation between adder bit-slices successfully demonstrates the capability to capture performance variation with great precision and relative simplicity. The measurements enabled by this circuitry show that decomposition of variation is often not straight-forward: circuit structure can introduce strong correlations that should not be associated with parameter variation. Similarly, systematic, position-dependent variation is separate from distance-dependent spatial correlation, but the two effects are not mutually exclusive.

When attempting to mitigate the impact of variation, the multi-dimensional nature of goals and constraints makes solutions complex. Performance, power/energy, area, design-complexity, and yield are all tightly coupled and must often be traded off amongst each other. Adaptive Voltage Scaling (AVS) and Adaptive Body-Biasing (ABB) both involve trade-offs of all of these parameters, particularly performance and power. Again, operating regime dictates the relative weights of each of the impacted

parameters. In the context of massively-parallel multicore systems, the simple inclusion of an additional power-supply voltage, together with the determination of the optimal value of that power-supply voltage, significantly reduces energy while guaranteeing performance and yield constraints, and minimizing area and design-complexity overheads.

## 6.2 Future Work

Though there has been great progress in understanding, characterizing and mitigating variation, much work remains as scaling continues. This concluding section provides some thoughts about specific projects that might be undertaken to improve understanding of variation and its impacts.

### 6.2.1 Devices

Recently, Cao et al. have produced a framework to create predictive technology models (PTM) for circuit designers to use before a new process exists, or if access to a process is not available [113]. These models typically include process corners, but do not include variation statistics in any form. Development of predictive variation models is crucial to the success of any predictive technology model as the relative magnitude of variation increases. In particular, analytic models based on modeling of fundamental physical processes or empirical modeling could greatly aid in predictive variation models for threshold voltage ($V_T$), channel length ($L$), oxide thickness ($t_{ox}$) and mobility ($\mu$).

This type of modeling has been done to some degree for threshold voltage by Asenov et al. and is shown in Eq. 2.8. Here, the nominal values of channel doping ($N_{ch}$), $t_{ox}$ and transistor dimensions provide an estimate of the variation in threshold voltage. These nominal values can be estimated from predictive technology models. Similar models should be developed, through fundamental physical modeling or extensive TCAD simulation, for channel length, oxide thickness and mobility variation. Ideally, all of these models would contain only process-determined variables. This

would also aid in process development, as process engineers could quickly identify the sensitivity of device parameters to variation in certain process parameters.

To achieve increased densities and device performance, processes, materials and device design continually change. In the near future, the industry is expected to significantly change the structure of transistors. The most common expectation is that a vertical channel, undoped device (either a FinFET or similar device) will replace the lateral channel MOSFET used for the past three to four decades. New device designs will have different sensitivities to various processing steps and the variation signatures of these new devices will differ as well. Characterization of the susceptibility of these new devices to variation as well as their variation signatures is critical in deciding which device structures to adopt. This has begun to a small extent in [51] but more work in this area is needed, especially for devices like carbon nanotube FETs (if manufacturable), carbon nanowires and other such devices, especially for interconnect applications.

## 6.2.2  Circuits

The following are specific areas of possible future research that follow from the data and analysis at the circuit level.

Standard cell libraries for ASIC flows using static timing analysis are often characterized by performing many Spice-level circuit simulations. Foundries that provide standard cell libraries for their process could make use of the extremely fine-resolution all-digital delay measurement technique described in Chapter 4 for this purpose. Although this would be limited by the cost of fabricating a test-die for this purpose, foundries that develop processes and run test-chips during process development could afford to do this. In addition to giving the foundry real data for both nominal timing and variation statistics, by characterizing different variants of standard cells, optimal standard cell libraries could be developed in parallel with the process. The technique would need modification to be able to measure rise/fall times, and additional circuitry would be needed to vary the rise/fall times of input signals, but the technique is otherwise directly applicable to such use.

To overcome some of the limitations being posed by deep sub-wavelength lithography, researchers at Carnegie Mellon University have proposed using "logic bricks" which are made up of larger blocks of common boolean logic rather than individual standard cells [130]. The layout of these bricks is highly optimized to achieve layout regularity and improve the lithographic printability of the layout. Spatial analysis of the variation in such layouts would be interesting to see if such regularity possibly introduces longer-range optical phenomena resulting in spatial correlation, or if the variation remains largely uncorrelated.

At the end of Chapter 4, adaptive body-biasing and voltage-scaling were evaluated as potential variation mitigation techniques. Although voltage-scaling is generally more effective, there may be some designs which may benefit from both techniques. In particular, highly duty-cycled systems which are in standby mode for large periods of time but require high throughput when active could be ideal candidates for joint optimization of adaptive body-biasing and voltage-scaling. In such situations, minimization of leakage is critical during periods of inactivity, but throughput and/or minimum retention voltages of memory elements may limit the magnitude of dynamic power-supply voltage scaling. These constraints make development of a framework that is capable of joint optimization of both techniques valuable.

To better perform any optimization where both leakage and active energy are involved, understanding the relative contributions and impact of variation on each is important. This could be achieved with a redesign of the second test-chip to include power/energy measurement capabilities for both idle and active components. The addition of simple power-gating circuitry would enable measurement of individual adder power through off-chip measurement. On-chip measurement circuitry would involve measuring the voltage drop across the power transistors, but is complicated by the need to accurately determine the resistance of those transistors. This simple redesign could be coupled with future work in the section below.

## 6.2.3 Architectures & Systems

At the architecture and systems level there is broad scope for future work, as the research community has only recently begun to address variation at this level. A few specific opportunities that directly follow from the work presented in Chapter 5 follow.

Instead of re-designing the second test-chip with only the addition of power measurement capabilities, the adders should be replaced with small processor cores (perhaps ARM or MIPS cores which are sufficiently simple to implement) and multiple voltage supplies. The larger number of potential critical paths would better represent variation data from a multi-core processor. The inclusion of multiple voltage supplies would allow for confirmation of the simulated energy savings presented in Chapter 5.

The custom simulation methodology presented in Section 5.3.2 and Figure 5.3.2 made use of a non-optimized standard cell library based on predictive technology models. This resulted in a sub-optimal netlist, as well as poor understanding of leakage versus active energy contributions. The design should be re-synthesized using transistor models and standard cell libraries from the most advanced commercially available processes with, ideally, complete implementation of SRAMs as well. The complete design should then be logically simulated using test vectors to characterize nodal activity factors for use in a power simulator. Additionally, a simulator such as HotLeakage [131] should be used to model leakage power, and to identify whether leakage can be safely ignored in the analytic framework in Section 5.2.1. If it cannot, the simulation results from HotLeakage and potentially measured data from the last item should be used as parameters for a formulation which includes leakage (see Section 5.3.4).

An outstanding issue in the analysis of multiple voltage supplies for energy savings with the granularity of individual cores is that there are perhaps only a few critical paths that require a higher voltage. By utilizing only core-level granularity, the energy savings is likely less than that of finer granularities [96]. However, fine granularity also potentially increases the amount of overhead, as power multiplexing

circuitry is now needed for each block (be it a gate, path or larger block). Studies are required to determine the optimal level of granularity which results in largest energy savings and smallest overhead. Examples of granularities that could be studied are individual cores, functional blocks, paths, just expected critical paths based on simulation results, and gates. The results of this investigation may also significantly alter the analytic formulation in Section 5.2.1.

Recently, architects have begun to develop tools that allow for architecture exploration, while considering process variation in addition to other metrics. Polaris, a tool to aid in exploration of on-chip interconnection networks, is one such tool [132]. A most-urgent and necessary addition to this tool is voltage-scaling as a means of addressing process variation — this will likely dramatically affect active power which is not currently addressed in the current implementation of the tool. Very simple models based on the equations found in Section 5.2.1 could be included in the tool.

# Appendix A

# Performance of MEVS Algorithm

## A.1 Proof of Optimality for $N = 2$

**Theorem A.1.1.** *There exists only one solution to the $N = 2$ case and that solution must lie in $0 \leq V_1 < V_2$.*

*Proof.* We begin by taking the derivative of Eq. 5.8 with respect to $V_1$ and setting it equal to 0, which gives:

$$\frac{f(V_1)}{F(V_1)} = \frac{2V_1}{V_2^2 - V_1^2} \tag{A.1}$$

The LHS of Eq. A.1 is shown to be a positive, strictly decreasing function by Pechtl in [133]. Furthermore, Pechtl also shows that it is asymptotic to $-V_1$ as $V_1 \rightarrow -\infty$ and goes to 0 as $V_1 \rightarrow \infty$. On the right side of Eq. A.1, the numerator is strictly increasing and the denominator is strictly decreasing over $0 \leq V_1 < V_2$, starting at $V_2^2$ when $V_1 = 0$ and reaching 0 at $V_1 = V_2$. Thus, the LHS is monotonically decreasing while the RHS is monotonically increasing, so there is at most one intersection point, and it must be located within the range of interest since the RHS ranges from 0 (at $V_1 = 0$) to $\infty$ (at $V_1 = V_2$). $\qquad \square$

## A.2  Bounds on MEVS Algorithm

To show the bounds on the overhead as defined in Section 5.2.2, we begin with the following. Let the $N$ voltages be distributed with intervals of $\delta_k := g(k)$. We define $g(x)$ to be the "continuous" form of $g(k)$ and $G(x) = \int_0^x g(y)dy$. We also define $V_k = V_L + \sum_{j=1}^{k} \delta_j$. Over any one interval, the energy cost in the discretized case is given by:

$$E_{int} = \int_{V_k}^{V_{k+1}} V_{k+1}^2 f(V)dV \tag{A.2}$$

Similarly, in the ideal, continuous case it is Eq. 5.12 over a single interval:

$$E_{int} = \int_{V_k}^{V_{k+1}} V^2 f(V)dV \tag{A.3}$$

Subtracting the two and using a change of variables ($V = V_k + t, => dv = dt$), we get:

$$E_{overhead,int} = \int_0^{\delta_k} \left[2V_k(\delta_k - t) + \delta_k^2 - t^2\right] f(V_k + t)dt \tag{A.4}$$

However, in the limit of large $N$ and small $\delta_k$, $\delta_k^2 - t^2$ is small enough to be ignored (this introduces an error on the order of $O(\frac{1}{N^2})$ which is smaller than the overall overhead, see below) and $f(V_k + t) \approx f(V_k)$. So, Eq. A.4 reduces to:

$$E_{overhead,int} \leq \delta_k^2 V_k f(V_k) \tag{A.5}$$

The total overhead is then the summation of the individual interval overheads:

$$E_{overhead,TOT} \leq \sum_{k=1}^{N} \delta_k^2 V_k f(V_k) \tag{A.6}$$

We choose $g(k)$ such that the size of each interval results in the voltage "cost", $V^2$, from Eq. 5.12 being balanced across all intervals. As will be seen, choosing $g(k)$ as

in Eq. A.7 will achieve this.

$$g(k) \propto \frac{1}{V_L + \sum_{j=1}^{k-1} g(j)} = \frac{1}{V_{k-1}} \tag{A.7}$$

In the continuous case, this becomes:

$$\frac{dG(x)}{dx} = \frac{C}{V_L + G(x)} \tag{A.8}$$

where $C$ is a normalization constant. To be self-consistent, we also require that $\sum_{k=1}^{N} g(k) = |V_U - V_L|$ in the discrete case and $G(N) = |V_U - V_L|$ in the continuous case.

Integrating Eq. A.8 we get:

$$G(x) = \sqrt{V_L^2 + 2Cx} - V_L \tag{A.9}$$

and using the boundary conditions, $C$ can be computed to be:

$$C = \frac{1}{2N} \left( V_U^2 - V_L^2 \right) \tag{A.10}$$

We can also make the following approximation:

$$\delta_k := g(k) \simeq G(k) - G(k-1) \simeq G'(k) \tag{A.11}$$

$$\simeq \frac{C}{V_L + G(k)} \tag{A.12}$$

which gives $C = G'(k) \left[ V_L + G(k) \right]$.

Since $\delta_k \simeq G'(k)$, we can reformulate Eq. A.6 as:

$$E_{overhead,TOT} = \sum_{k=1}^{N} G'(k)^2 \left( V_L + G(k) \right) f \left( V_L + G(k) \right) \tag{A.13}$$

153

and using Eq. A.11 this becomes:

$$E_{overhead,TOT} = C \left[ \sum_{k=1}^{N} G'(k) f\left(V_L + G(k)\right) \right] \qquad (A.14)$$

the latter half of which, in the limit of large $N$, is a Riemann integral:

$$E_{overhead,TOT} \simeq C \int_{V_L}^{V_U} f(V) dV \qquad (A.15)$$

Since the integral in the above equation is simply $F(V_U) - F(V_L) < 1$, we can say:

$$E_{overhead,TOT} \leq C = \frac{1}{2N}(V_U^2 - V_L^2) \qquad (A.16)$$

$$\leq \frac{1}{2N}(V_U - V_L)(V_U + V_L) \qquad (A.17)$$

Finally, in general $V_U - V_L \propto \sigma$ and $V_U + V_L \propto \mu$, so $E_{overhead,TOT} \leq O(\frac{\mu\sigma}{2N})$.

# Bibliography

[1] S. Borkar, T. Karnik, S. Narendra, J. Tschanz, A. Keshavarzi, and V. De, "Parameter Variations and Impact on Circuits and Microarchitecture," pp. 338–342, Proceedings of the Design Automation Conference, 2003.

[2] M. Quirk and J. Serda, *Semiconductor Manufacturing Technology*. Prentice Hall, 2001.

[3] A. Asenov, "Statistical Device Variability and Its Impact on Design," International Symposium on Asynchronous Circuits and Systems, April 2008.

[4] K. Kuhn, C. Kenyon, A. Kornfeld, M. Liu, A. Maheshwari, W.-k. Shih, S. Sivakumar, G. Taylor, P. VanDerVoorn, and K. Zawadzki, "Managing Process Variation in Intels 45nm CMOS Technology," *Intel Technology Journal*, vol. 12, June 2008.

[5] J. del Alamo, *Integrated Microelectronic Devices*. Prentice Hall, 2007.

[6] "International Technology Roadmap for Semiconductors 2001 Edition," tech. rep., International Technology Roadmap for Semiconductors, 2001.

[7] "International Technology Roadmap for Semiconductors 2003 Edition," tech. rep., International Technology Roadmap for Semiconductors, 2003.

[8] "International Technology Roadmap for Semiconductors 2005 Edition," tech. rep., International Technology Roadmap for Semiconductors, 2005.

[9] "International Technology Roadmap for Semiconductors 2007 Edition," tech. rep., International Technology Roadmap for Semiconductors, 2007.

[10] K. Kuhn, "Reducing Variation in Advanced Logic Technologies: Approaches to Process and Design for Manufacturability of Nanoscale CMOS," pp. 471–474, IEEE International Electron Devices Meeting, Dec. 2007.

[11] F. Hamzaoglu, K. Zhang, Y. Wang, H. J. Ahn, U. Bhattacharya, Z. Chen, Y.-G. Ng, A. Pavlov, K. Smits, and M. Bohr, "A 3.8 GHz 153 Mb SRAM Design With Dynamic Stability Enhancement and Leakage Reduction in 45 nm High-k Metal Gate CMOS Technology," *IEEE Journal of Solid-State Circuits*, vol. 44, pp. 148–154, January 2009.

[12] S. Das, C. Tokunaga, S. Pant, W.-H. Ma, S. Kalaiselvan, K. Lai, D. M. Bull, and D. T. Blaauw, "RazorII: In Situ Error Detection and Correction for PVT and SER Tolerance," *IEEE Journal of Solid-State Circuits*, vol. 44, pp. 32–48, January 2009.

[13] K. Bernstein, D. J. Frank, A. E. Gattiker, W. Haensch, B. L. Ji, S. R. Nassif, E. J. Nowak, D. J. Pearson, and N. J. Rohrer, "High-Performance CMOS Variability in the $65 - nm$ Regime and Beyond," *IBM Journal of Research and Development*, vol. 50, no. 4/5, pp. 433–449, 2006.

[14] G. E. Moore, "Cramming More Components Onto Integrated Circuits," *Electronics*, vol. 38, April 1965.

[15] S. Borkar, "Designing Reliable Systems from Unreliable Components: The Challenges of Transistor Variability and Degradation," *IEEE Micro*, vol. 25, pp. 10–16, November-December 2005.

[16] P. Friedberg, *Spatial Modeling of Gate Length Variation for Process-Design Co-Optimization*. PhD dissertation, Univ. of California, Berkeley, Dept. of Electrical Engineering and Computer Science, 2007.

[17] S. R. Nassif, "Design for Variability in DSM Technologies," pp. 451–454, IEEE International Symposium on Quality Electronic Design, 2000.

[18] X.-W. Lin, "Design and Process Variability - the Sources and Mechanisms," DAC Tutotrial - Practical Aspects of Coping with Variability: An Electrical View, 2006.

[19] L. van den Hove, K. Ronse, and R. Pforr, "Optical Lithography Techniques for $0.25\mu m$ and Below: CD Control Issues," pp. 24–30, International Symposium on VLSI Technology, Systems, and Applications, May 1995.

[20] P. Gupta, A. Kahng, Y. Kim, and D. Sylvester, "Self-Compensating Design for Focus Variation," pp. 365–368, Design Automation Conference, June 2005.

[21] L. W. Liebmann, S. M. Mansfield, A. K. Wong, M. A. Lavin, W. C. Leipold, and T. G. Dunham, "TCAD Development for Lithography Resolution Enhancement," *IBM Journal of Research and Development*, vol. 45, pp. 651–665, September 2001.

[22] A. Asenov, S. Kaya, and A. Brown, "Intrinsic Parameter Fluctuations in Decananometer MOSFETs Introduced by Gate Line Edge Roughness," *IEEE Transactions on Electron Devices*, vol. 50, pp. 1254–1260, May 2003.

[23] H. Fukuda, "Analysis of Line Edge Roughness Using Probability Process Model for Chemically Amplified Resists," pp. 76–77, International Microprocesses and Nanotechnology Conference, November 2002.

[24] G. May, J. Huang, and C. Spanos, "Statistical Experimental Design in Plasma Etch Modeling," *IEEE Transactions on Semiconductor Manufacturing*, vol. 4, pp. 83–98, May 1991.

[25] K. Abrokwah, "Characterization and Modeling of Plasma Etch Pattern Dependencies in Integrated Circuits," MEng thesis, Massachusetts Institute of Technology, Dept. of Elect. Engineering and Comp. Science, June 2006.

[26] R. W. Keyes, "The Effect of Randomness in the Distribution of Impurity Atoms on FET Thresholds," *Applied Physics*, vol. 8, pp. 251–259, 1975.

[27] X. Xie, D. Boning, F. Meyer, and R. Rzehak, "Analysis of Nanotopography and Layout Variations in Patterned STI CMP," International Conference on Planarization/CMP Technology, 2006.

[28] X. Xie, *Physical Understanding and Modeling of Chemical Mechanical Planarization in Dielectric Materials*. PhD dissertation, Massachusetts Institute of Technology, Dept. of Physics, June 2007.

[29] D. Boning, W. Moyne, T. Smith, J. Moyne, R. Telfeyan, A. Hurwitz, S. Shellman, and J. Tayor, "Run by Run Control of Chemical-Mechanical Polishing," *IEEE Transactions on Components, Packaging, and Manufacturing Technology, Part C*, vol. 19, pp. 307–314, October 1996.

[30] B. Stine, D. Boning, J. Chung, L. Camilletti, F. Kruppa, E. Equi, W. Loh, S. Prasad, M. Muthukrishnan, D. Towery, M. Berman, and A. Kapoor, "The Physical and Electrical Effects of Metal-Fill Patterning Practices for Oxide Chemical-Mechanical Polishing Processes," *IEEE Transactions on Electron Devices*, vol. 45, pp. 665–679, March 1998.

[31] R. Deaton and H. Massoud, "Manufacturability of Rapid-Thermal Oxidation of Silicon: Oxide Thickness, Oxide Thickness Variation, and System Dependency," *IEEE Transactions on Semiconductor Manufacturing*, vol. 5, pp. 347–358, November 1992.

[32] J. Hoyt, H. Nayfeh, S. Eguchi, I. Aberg, G. Xia, T. Drake, E. Fitzgerald, and D. Antoniadis, "Strained Silicon MOSFET Technology," pp. 23–26, IEEE International Electron Devices Meeting, 2002.

[33] P. Bai, C. Auth, S. Balakrishnan, M. Bost, R. Brain, V. Chikarmane, R. Heussner, M. Hussein, J. Hwang, D. Ingerly, R. James, J. Jeong, C. Kenyon, E. Lee, S.-H. Lee, N. Lindert, M. Liu, Z. Ma, T. Marieb, A. Murthy, R. Nagisetty, S. Natarajan, J. Neirynck, A. Ott, C. Parker, J. Sebastian, R. Shaheed, S. Sivakumar, J. Steigerwald, S. Tyagi, C. Weber, B. Woolery, A. Yeoh, K. Zhang, and M. Bohr, "A 65$nm$ Logic Technology Featuring 35nm Gate Lengths, Enhanced Channel Strain, 8 Cu Interconnect Layers, low-k ILD and 0.57 $\mu m^2$ SRAM Cell," pp. 657–660, IEEE International Electron Devices Meeting, December 2004.

[34] C. Gallon, G. Reimbold, G. Ghibaudo, R. Blanchi, R. Gwoziecki, and C. Raynaud, "Electrical Analysis of Mechanical Stress Induced by Shallow Trench Isolation MOSFETs," pp. 359–362, European Solid-State Device Research Conference, September 2003.

[35] W. Zhao, F. Liu, K. Agarwal, D. Acharyya, S. Nassif, K. Nowka, and Y. Cao, "Rigorous Extraction of Process Variations for 65-nm CMOS Design," *IEEE Transactions on Semiconductor Manufacturing*, vol. 22, pp. 196–203, February 2009.

[36] J. Chen, T.-Y. Chan, P. Ko, and C. Hu, "Gate Current in OFF-state MOSFET," *IEEE Electron Device Letters*, vol. 10, pp. 203–205, May 1989.

[37] S. Datta, G. Dewey, M. Doczy, B. Doyle, B. Jin, J. Kavalieros, R. Kotlyar, M. Metz, N. Zelick, and R. Chau, "High Mobility Si/SiGe Strained Channel MOS Transistors with HfO2/TiN Gate Stack," pp. 28.1.1–28.1.4, IEEE International Electron Devices Meeting, December 2003.

[38] K. Mistry, C. Allen, C. Auth, B. Beattie, D. Bergstrom, M. Bost, M. Brazier, M. Buehler, A. Cappellani, R. Chau, C.-H. Choi, G. Ding, K. Fischer, T. Ghani, R. Grover, W. Han, D. Hanken, M. Hattendorf, J. He, J. Hicks, R. Huessner, D. Ingerly, P. Jain, R. James, L. Jong, S. Joshi, C. Kenyon, K. Kuhn, K. Lee, H. Liu, J. Maiz, B. McIntyre, P. Moon, J. Neirynck, S. Pae, C. Parker, D. Parsons, C. Prasad, L. Pipes, M. Prince, P. Ranade, T. Reynolds, J. Sandford, L. Shifren, J. Sebastian, J. Seiple, D. Simon, S. Sivakumar, P. Smith, C. Thomas, T. Troeger, P. Vandervoorn, S. Williams, and K. Zawadzki, "A 45nm Logic Technology with High-k+Metal Gate Transistors, Strained Silicon, 9 Cu Interconnect Layers, 193nm Dry Patterning, and 100% Pb-free Packaging," pp. 247–250, IEEE International Electron Devices Meeting, December 2007.

[39] M. Cho, K. Maitra, and S. Mukhopadhyay, "Analysis of the Impact of Interfacial Oxide Thickness Variation on Metal-Gate High-K Circuits," pp. 285–288, IEEE Custom Integrated Circuits Conference, September 2008.

[40] R. Sokel, "Transistor Scaling with Constant Subthreshold Leakage," *IEEE Electron Device Letters*, vol. 4, pp. 85–87, April 1983.

[41] M. J. M. Pelgrom, A. C. J. Duinmaijer, and A. P. G. Welbers, "Matching Properties of MOS Transistors," *IEEE Journal of Solid-State Circuits*, vol. 24, pp. 1433–1440, October 1989.

[42] A. Asenov, A. Brown, J. Davies, S. Kaya, and G. Slavcheva, "Simulation of Intrinsic Parameter Fluctuations in Decananometer and Nanometer-Scale MOSFETs," *IEEE Transactions on Electron Devices*, vol. 50, pp. 1837–1852, September 2003.

[43] B. E. Deal, M. Slkar, A. S. Grove, and E. H. Snow, "Characteristics of the Surface-State Charge $(Q_{ss})$ of Thermally Oxidized Silicon," *Journal of The Electrochemical Society*, vol. 114, pp. 266–274, March 1967.

[44] A. Gu and A. Zakhor, "Optical Proximity Correction With Linear Regression," *IEEE Transactions on Semiconductor Manufacturing*, vol. 21, pp. 263–271, May 2008.

[45] H. Cai, *Modeling of Pattern Dependencies in the Fabrication of Multilevel Copper Metallization*. PhD dissertation, Massachusetts Institute of Technology, Dept. of Materials Science and Engineering, June 2007.

[46] R. Berry, "Correlation of Diffusion Process Variations with Variations in Electrical Parameters of Bipolar Transistors," *Proceedings of the IEEE*, vol. 57, pp. 1513–1517, September 1969.

[47] M. Pocha, A. Gonzalez, and R. Dutton, "Threshold Voltage Controllability in Double-Diffused-MOS Transistors," *IEEE Transactions on Electron Devices*, vol. 21, pp. 778–784, December 1974.

[48] W. Schemmert and G. Zimmer, "Threshold-Voltage Sensitivity of Ion-Implanted M.O.S. Transistors Due to Process Variations," *Electronics Letters*, vol. 10, pp. 151–152, February 1974.

[49] S. Lin and C. Wong, "Process-Variation-Tolerant Zero Skew Clock Routing," pp. 83–86, IEEE International ASIC Conference and Exhibit, September 1993.

[50] S. Nassif, "Modeling and Forecasting of Manufacturing Variations," pp. 2–10, International Workshop on Statistical Metrology, 2000.

[51] S. Xiong and J. Bokor, "Sensitivity of Double-Gate and FinFET Devices to Process Variations," *IEEE Transactions on Electron Devices*, vol. 50, pp. 2255–2261, November 2003.

[52] A. Srivastava, D. Sylvester, and D. Blaauw, *Statistical Analysis and Optimization for VLSI: Timing and Power*. New York: Springer, 2007.

[53] K. Gettings and D. Boning, "Study of CMOS Process Variation by Multiplexing Analog Characteristics," *IEEE Transactions on Semiconductor Manufacturing*, vol. 21, pp. 513–525, November 2008.

[54] K. Agarwal, F. Liu, C. McDowell, S. Nassif, K. Nowka, M. Palmer, D. Acharyya, and J. Plusquellic, "A Test Structure for Characterizing Local Device Mismatches," pp. 67–68, Symposium on VLSI Circuits, 2006.

[55] K. Agarwal, J. Hayes, and S. Nassif, "Fast Characterization of Threshold Voltage Fluctuation in MOS Devices," *IEEE Transactions on Semiconductor Manufacturing*, vol. 21, pp. 526–533, November 2008.

159

[56] R. Rao, K. Jenkins, and J.-J. Kim, "A Completely Digital On-Chip Circuit for Local-Random-Variability Measurement," pp. 412–623, IEEE International Solid-State Circuits Conference, February 2008.

[57] S. Mukhopadhyay, K. Kim, K. Jenkins, C.-T. Chuang, and K. Roy, "Statistical Characterization and On-Chip Measurement Methods for Local Random Variability of a Process Using Sense-Amplifier-Based Test Structure," pp. 400–611, IEEE International Solid-State Circuits Conference, February 2007.

[58] T. Fischer, E. Amirante, P. Huber, T. Nirschl, A. Olbrich, M. Ostermayr, and D. Schmitt-Landsiedel, "Analysis of Read Current and Write Trip Voltage Variability From a 1-MB SRAM Test Structure," *IEEE Transactions on Semiconductor Manufacturing*, vol. 21, pp. 534–541, November 2008.

[59] V. Wang and K. Shepard, "On-Chip Transistor Characterisation Arrays for Variability Analysis," *Electronics Letters*, vol. 43, no. 15, pp. 806–807, 2007.

[60] J. S. Panganiban, "A Ring Oscillator Based Variation Test Chip," MEng thesis, Massachusetts Institute of Technology, Dept. of Elect. Engineering and Comp. Science, June 2002.

[61] L.-T. Pang and B. Nikolic, "Impact of Layout on 90nm CMOS Process Parameter Fluctuations," pp. 69–70, Symposium on VLSI Circuits, 2006.

[62] L.-T. Pang and B. Nikolic, "Measurement and Analysis of Variability in 45nm Strained-Si CMOS Technology," pp. 129–132, IEEE Custom Integrated Circuits Conference, September 2008.

[63] M. Bhushan, M. Ketchen, S. Polonsky, and A. Gattiker, "Ring Oscillator Based Technique for Measuring Variability Statistics," pp. 87–92, IEEE International Conference on Microelectronic Test Structures, March 2006.

[64] Z. Guo, A. Carlson, L.-T. Pang, K. Duong, T.-J. K. Liu, and B. Nikolic, "Large-Scale Read/Write Margin Measurement in 45nm CMOS SRAM Arrays," pp. 42–43, IEEE Symposium on VLSI Circuits, June 2008.

[65] H.-Y. Wong, L. Cheng, Y. Lin, and L. He, "FPGA Device and Architecture Evaluation Considering Process Variations," pp. 19–24, IEEE/ACM International Conference on Computer-Aided Design, November 2005.

[66] M. Annavaram, E. Grochowski, and P. Reed, "Implications of Device Timing Variability on Full Chip Timing," pp. 37–45, International Symposium on High Performance Computer Architecture, 2007.

[67] E. Humenay, D. Tarjan, and K. Skadron, "Impact of Process Variations on Multicore Performance Symmetry," pp. 1–6, Design, Automation & Test in Europe Conference & Exhibition, 2007.

[68] K. Meng, F. Huebbers, R. Joseph, and Y. Ismail, "Modeling and Characterizing Power Variability in Multicore Architectures," pp. 146–157, IEEE International Symposium on Performance Analysis of Systems & Software, April 2007.

[69] J. Kibarian and A. Strojwas, "Using Spatial Information to Analyze Correlations Between Test Structure Data," pp. 187–191, International Conference on Microelectronic Test Structures, March 1990.

[70] P. Friedberg, Y. Cao, J. Cain, R. Wang, J. Rabaey, and C. Spanos, "Modeling Within-Die Spatial Correlation Effects for Process-Design Co-optimization," pp. 516–521, International Symposium on Quality of Electronic Design, March 2005.

[71] C. Conroy, W. Lane, and M. Moran, "Statistical Design Techniques for D/A Converters," *IEEE Journal of Solid-State Circuits*, vol. 24, pp. 1118–1128, August 1989.

[72] R. Guldi, C. Jenkins, G. Oamminga, T. Baum, and T. Foster, "Process Optimization Tweaking Tool (POTT) and its Application in Controlling Oxidation Thickness," *IEEE Transactions on Semiconductor Manufacturing*, vol. 2, pp. 54–59, May 1989.

[73] S. Leang and C. Spanos, "Statistically Based Feedback Control of Photoresist Application," pp. 185–190, IEEE/SEMI Advanced Semiconductor Manufacturing Conference and Workshop, October 1991.

[74] K. S. Wong, D. S. Boning, and H. H. Sawin, "On Endpoint Detection of Plasma Etching via Optical Emission Interferometry," p. 302, Electrochemical Society Meeting, May 1996.

[75] S. Owa, K. Nakano, H. Nagasaka, T. Fujiwara, T. Matsuyama, Y. Ohmura, and H. Magoona, "Immersion Lithography Ready for $45nm$ Manufacturing and Beyond," pp. 238–244, IEEE/SEMI Advanced Semiconductor Manufacturing Conference, June 2007.

[76] N. Verma and A. Chandrakasan, "A 256 kb 65 nm 8T Subthreshold SRAM Employing Sense-Amplifier Redundancy," *IEEE Journal of Solid-State Circuits*, vol. 43, pp. 141–149, January 2008.

[77] N. Verma and A. P. Chandrakasan, "A High-Density 45 nm SRAM Using Small-Signal Non-Strobed Regenerative Sensing," *IEEE Journal of Solid-State Circuits*, vol. 44, pp. 163–173, January 2009.

[78] J. Tschanz, J. Kao, S. Narendra, R. Nair, D. Antoniadis, A. Chandrakasan, and V. De, "Adaptive body bias for reducing impacts of die-to-die and within-die parameter variations on microprocessor frequency and leakage," *IEEE Journal of Solid-State Circuits*, vol. 37, pp. 1396–1402, November 2002.

[79] T. Burd, T. Pering, A. Stratakos, and R. Brodersen, "A Dynamic Voltage Scaled Microprocessor System," *IEEE Journal of Solid-State Circuits*, vol. 35, pp. 1571–1580, November 2000.

[80] R. Gonzalez, B. Gordon, and M. Horowitz, "Supply and Threshold Voltage Scaling for Low Power CMOS," *IEEE Journal of Solid-State Circuits*, vol. 32, pp. 1210–1216, August 1997.

[81] M. Eireiner, S. Henzler, G. Georgakos, J. Berthold, and D. Schmitt-Landsiedel, "In-Situ Delay Characterization and Local Supply Voltage Adjustment for Compensation of Local Parametric Variations," *IEEE Journal of Solid-State Circuits*, vol. 42, pp. 1583–1592, July 2007.

[82] A. Agarwal, B. Paul, S. Mukhopadhyay, and K. Roy, "Process Variation in Embedded Memories: Failure Analysis and Variation Aware Architecture," *IEEE Journal of Solid-State Circuits*, vol. 40, pp. 1804–1814, September 2005.

[83] E. Fetzer, "Using Adaptive Circuits to Mitigate Process Variations in a Microprocessor Design," *IEEE Design & Test of Computers*, vol. 23, pp. 476–483, June 2006.

[84] D. Marculescu and S. Garg, "Process-Driven Variability Analysis of Single and Multiple VoltageFrequency Island Latency-Constrained Systems," *IEEE Transactions on Computer Aided Design*, vol. 27, pp. 893–905, May 2008.

[85] P. Mozumder, C. Shyamsundar, and A. Strojwas, "Statistical control of VLSI fabrication processes. A framework," *IEEE Transactions on Semiconductor Manufacturing*, vol. 1, pp. 62–71, May 1988.

[86] A. Gregene and H. Camenzind, "Frequency-Selective Integrated Circuits Using Phase-Lock Techniques," *IEEE Journal of Solid-State Circuits*, vol. 4, pp. 216–225, August 1969.

[87] D. Harris and S. Naffziger, "Statistical Clock Skew Modeling with Data Delay Variations," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 9, pp. 888–898, December 2001.

[88] N. Drego, "A Low-Skew, Low-Jitter Receiver Circuit for On-Chip Optical Clock Distribution," Masters thesis, Massachusetts Institute of Technology, Dept. of Elect. Engineering and Comp. Science, June 2003.

[89] M. J. Kobrinsky, B. A. Block, J.-F. Zheng, B. C. Barnett, E. Mohammed, M. Reshotko, F. Robertson, S. List, I. Young, and K. Cadien, "On-Chip Optical Interconnects," *Intel Technology Journal*, vol. 8, pp. 129–141, May 2004.

[90] B. Calhoun and A. Chandrakasan, "A 256kb Sub-threshold SRAM in 65nm CMOS," pp. 2592–2601, IEEE International Solid-State Circuits Conference,, February 2006.

[91] S. Das, D. Roberts, S. Lee, S. Pant, D. Blaauw, T. Austin, K. Flautner, and T. Mudge, "Self-Tuning DVS Processor Using Delay-Error Detection and Correction," *IEEE Journal of Solid-State Circuits*, vol. 41, pp. 792–804, April 2006.

[92] N. Verma, J. Kwong, and A. Chandrakasan, "Nanometer MOSFET Variation in Minimum Energy Subthreshold Circuits," *IEEE Transactions on Electron Devices*, vol. 55, pp. 163–174, January 2008.

[93] S. Pandey, R. Drechsler, T. Murgan, and M. Glesner, "Process Variations Aware Robust On-Chip Bus Architecture Synthesis for MPSoCs," pp. 2989–2992, IEEE International Symposium on Circuits and Systems, May 2008.

[94] G. Nabaa, N. Azizi, and F. Najm, "An Adaptive FPGA Architecture with Process Variation Compensation and Reduced Leakage," pp. 624–629, ACM/IEEE Design Automation Conference, 2006.

[95] J. Donald and M. Martonosi, "Power Efficiency for Variation-Tolerant Multicore Processors," pp. 304–309, International Symposium on Low-Power Electronic Design, 2006.

[96] X. Liang, G.-Y. Wei, and D. Brooks, "ReVIVaL: A Variation-Tolerant Architecture Using Voltage Interpolation and Variable Latency," *SIGARCH Comput. Archit. News*, vol. 36, no. 3, pp. 191–202, 2008.

[97] J. Xiong, V. Zolotov, and L. He, "Robust Extraction of Spatial Correlation," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 26, no. 4, pp. 619–631, 2007.

[98] A. Agarwal, D. Blaauw, V. Zolotov, S. Sundareswaran, M. Zhao, K. Gala, and R. Panda, "Statistical Delay Computation Considering Spatial Correlations," pp. 271–276, Asia and South Pacific Design Automation Conference, 2003.

[99] H. Chang and S. S. Sapatnekar, "Full-Chip Analysis of Leakage Power Under Process Variations, Including Spatial Correlations," (New York, NY, USA), pp. 523–528, ACM, 2005.

[100] H. Chang and S. S. Sapatnekar, "Statistical Timing Analysis Under Spatial Correlations," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 24, no. 9, pp. 1467–1482, 2005.

[101] J. Fan, N. Mi, S. X. D. Tan, Y. Cai, and X. Hong, "Statistical Model Order Reduction for Interconnect Circuits Considering Spatial Correlations," pp. 1–6, Design, Automation & Test in Europe Conference & Exhibition, 2007.

[102] L. Zhang, Y. Hu, and C.-P. Chen, "Statistical Timing Analysis with Path Reconvergence and Spatial Correlations," p. 5, Design, Automation & Test in Europe Conference & Exhibition, 2006.

[103] F. Liu, "A General Framework for Spatial Correlation Modeling in VLSI Design," pp. 817–822, Design Automation Conference, 2007.

[104] H. Li, C.-K. Koh, V. Balakrishnan, and Y. Chen, "Statistical Timing Analysis Considering Spatial Correlations," pp. 102–107, International Symposium on Quality Electronic Design, 2007.

[105] A. Agarwal, D. Blaauw, and V. Zolotov, "Statistical Timing Analysis for Intra-Die Process Variations with Spatial Correlations," pp. 900–907, International Conference on Computer Aided Design, 2003.

[106] Synopsis, *HSPICE Simulation and Analysis User Guide*, March 2007.

[107] Z.-H. Liu, C. Hu, J.-H. Huang, T.-Y. Chan, M.-C. Jeng, P. Ko, and Y. Cheng, "Threshold Voltage Model for Deep-Submicrometer MOSFETs," *IEEE Transactions on Electron Devices*, vol. 40, pp. 86–95, January 1993.

[108] D. A. Johns and K. Martin, *Analog Integrated Circuit Design*. New York: John Wiley and Sons, Inc., 1997.

[109] J. Kwong and A. Chandrakasan, "Variation-Driven Device Sizing for Minimum Energy Subthreshold Circuits," pp. 8–13, International Symposium on Low-Power Electronic Design, 2006.

[110] B. Zhai, S. Hanson, D. Blaauw, and D. Sylvester, "Analysis and Mitigation of Variability in Subthreshold Design," pp. 20–25, International Symposium on Low-Power Electronic Design, 2005.

[111] B. Zhai, D. Blaauw, D. Sylvester, and K. Flautner, "The Limit of Dynamic Voltage Scaling and Insomniac Dynamic Voltage Scaling," *IEEE Transactions on VLSI Systems*, vol. 13, pp. 1239–1252, November 2005.

[112] B. Calhoun and A. Chandraksan, "Ultra-Dynamic Voltage Scaling (UDVS) Using Sub-Threshold Operation and Local Voltage Dithering," *IEEE Journal of Solid-State Circuits*, vol. 41, pp. 238–245, January 2006.

[113] Y. Cao, T. Sato, D. Sylvester, M. Orshansky, and C. Hu, "New Paradigm of Predictive MOSFET and Interconnect Modeling for Early Circuit Design," pp. 201–204, Custom Integrated Circuits Conference, 2000.

[114] D. Blaauw, S. Kalaiselvan, K. Lai, W.-H. Ma, S. Pant, C. Tokunaga, S. Das, and D. Bull, "Razor II: In Situ Error Detection and Correction for PVT and SER Tolerance," pp. 400–401, IEEE International Solid-State Circuits Conference, February 2008.

[115] S. Maggioni, A. Veggetti, A. Bogliolo, and L. Croce, "Random Sampling for On-Chip Characterization of Standard-Cell Propagation Delay," pp. 41–45, International Symposium on Quality Electronic Design, March 2003.

[116] R. Bhatti, M. Denneau, and J. Draper, "Duty Cycle Measurement and Correction Using a Random Sampling Technique," IEEE Midwest Symposium on Circuits and Systems, August 2005.

[117] R. Bhatti, M. Denneau, and J. Draper, "Phase Measurement and Adjustment of Digital Signals Using Random Sampling Technique," IEEE International Symposium on Circuits And Systems, May 2006.

[118] W.-H. Lee, J.-D. Cho, and S.-D. Lee, "A High Speed and Low Power Phase-Frequency Detector and Charge-Pump," pp. 269–272, Asia and South Pacific Design Automation Conference, January 1999.

[119] N. Drego, A. Chandrakasan, and D. Boning, "Lack of Spatial Correlation in MOSFET Threshold-Voltage Variation and Implications for Voltage Scaling," *To appear in IEEE Transactions on Semiconductor Manufacturing*, May 2009.

[120] B. Carlson and B. Giolma, "SmartReflex Power and Performance Management Technologies: reduced power consumption, optimized performance," tech. rep., Texas Instruments, 2008.

[121] D. Wentzlaff and A. Agarwal, "The Case for a Factored Operating System (FOS)," tech. rep., MIT Computer Science and Artificial Intelligence Laboratory (CSAIL), 2008.

[122] R. Manohar, "Reconfigurable Asynchronous Logic," pp. 13–20, IEEE Custom Integrated Circuits Conference, September 2006.

[123] K. Bowman, X. Tang, J. Eble, and J. Meindl, "Impact of Extrinsic and Intrinsic Parameter Variations on CMOS System on a Chip Performance," pp. 267–271, IEEE International ASIC/SOC Conference, 1999.

[124] S. Borkar, "Thousand Core Chips: A Technology Perspective," pp. 746–749, Proceedings of the Design Automation Conference, 2007.

[125] M. Taylor, J. Psota, A. Saraf, N. Shnidman, V. Strumpen, M. Frank, S. Amarasinghe, A. Agarwal, W. Lee, J. Miller, D. Wentzlaff, I. Bratt, B. Greenwald, H. Hoffmann, P. Johnson, and J. Kim, "Evaluation of the RAW Microprocessor: An Exposed-Wire-Delay Architecture for ILP and Streams," pp. 2–13, International Symposium on Computer Architecture, June 2004.

[126] J. Stine, I. Castellanos, M. Wood, J. Henson, F. Love, W. Davis, P. Franzon, M. Bucher, S. Basavarajaiah, J. Oh, and R. Jenkal, "FreePDK: An Open-Source Variation-Aware Design Kit," IEEE International Conference on Microelectronic Systems Education, 2007.

[127] S. Gunther and R. Singhal, "Next Generation Intel Microarchitecture (Nehalem) Family: Architectural Insights and Power Management," Intel Developer Forum, August 2008.

[128] J. Chang, M. Huang, J. Shoemaker, J. Benoit, S.-L. Chen, W. Chen, S. Chiu, R. Ganesan, G. Leong, V. Lukka, S. Rusu, and D. Srivastava, "The 65-nm 16-MB Shared On-Die L3 Cache for the Dual-Core Intel Xeon Processor 7100 Series," *IEEE Journal of Solid-State Circuits*, vol. 42, pp. 846–852, April 2007.

[129] W. Kim, M. Gupta, G.-Y. Wei, and D. Brooks, "System Level Analysis of Fast, Per-Core DVFS Using On-Chip Switching Regulators," pp. 123–134, IEEE International Symposium on High Performance Computer Architecture, February 2008.

[130] K. Y. Tong, V. Rovner, L. Pileggi, and V. Kheterpal, "Design Methodology of Regular Logic Bricks for Robust Integrated Circuits," pp. 162–167, International Conference on Computer Design, October 2006.

[131] Y. Zhang, D. Parikh, K. Sankaranarayanan, K. Skadron, and M. Stan, "HotLeakage: A Temperature-Aware Model of Subthreshold and Gate Leakage for Architects," tech. rep., University of Virginia, Dept. of Computer Science, 2003.

[132] V. Soteriou, N. Eisley, H. Wang, B. Li, and L.-S. Peh, "Polaris: A System-Level Roadmap for On-Chip Interconnection Networks," IEEE International Conference on Computer Design, 2006.

[133] A. Pechtl, "A Note on the Derivative of the Normal Distribution's Logarithm," *Archiv der Mathematik*, vol. 70, pp. 83–88, January 1998.