# AN ABSTRACT OF THE DISSERTATION OF

Robert Pawlowski for the degree of Doctor of Philosophy in Electrical and Computer Engineering presented on July 14 2014.

Title: Measurement and Analysis of Soft Error Vulnerability of Low-Voltage Logic and Memory Circuits

Abstract approved: _____

Patrick Yin Chiang

Scaling the supply voltage into the sub/near-threshold domain is one of the most effective methods for improving the energy efficiency of next-generation electronic microsystems. Unfortunately, the relationship between low-voltage operation and radiation-induced soft error rate is not widely known, as little research has been previously performed and reported for soft-error susceptibility of on-chip memory and logic at very low supply voltages. This information is critical for low-voltage circuit designers, as many applications that would benefit from the energy efficiency of sub/near-threshold also require high reliability. This work first details the design and implementation of a portable soft error reference platform, specifically targeting very low-voltage operation. The circuit-level details of a TSMC 65nm test-chip design are given, along with an analysis of data from experiments performed at Los Alamos Neutron Science Center (LANSCE) and the OSU radi-

ation center. Once this soft-error rate is known, error resiliency techniques must be utilized for increased processor reliability. The design and implementation of an error-resilient, near-threshold SIMD processor in an IBM 45nm SOI process will also be covered. This prototype demonstrates both increased reliability and improved throughput over a conventional SIMD pipeline while operating in near-threshold.

Measurement and Analysis of Soft Error Vulnerability of
Low-Voltage Logic and Memory Circuits

by

Robert Pawlowski

A DISSERTATION

submitted to

Oregon State University

in partial fulfillment of
the requirements for the
degree of

Doctor of Philosophy

Presented July 14 2014
Commencement June 2015

Doctor of Philosophy dissertation of Robert Pawlowski presented on
July 14 2014.

APPROVED:

_____

Major Professor, representing Electrical and Computer Engineering


_____

Director of the School of Electrical Engineering and Computer Science


_____

Dean of the Graduate School




I understand that my dissertation will become part of the permanent collection
of Oregon State University libraries. My signature below authorizes release of my
dissertation to any reader upon request.


_____

Robert Pawlowski, Author

# ACKNOWLEDGEMENTS

Also, to all of the friends I have made throughout graduate school, thank you for all of the time spent playing soccer, going to conferences, going out for food/drinks, etc. You have all made it an experience I will always cherish.

Thank you to my fiancee, Laurel Jones, for putting up with my late nights and times of high stress. I can't imagine making it through the past 5 years without you, and I'm excited for the next step in our lives together.

Cynthia, Bryan, and Danny. Thanks for being great siblings. You all have been supportive and fun, and I love you guys very much. Thank you to all of my grandparents for your support, as well.

Finally, to my parents: Thank you for your love and help throughout my entire life. You both have believed in me and provided words of guidance for me whenever I've needed it. This work would not have been possible without your support!

# CONTRIBUTION OF AUTHORS

Dr. Joseph Crop was instrumental to the success of the work presented in Chapter 2.

Dr. Evgeni Krimer, Dr. Nariman Moezzi-Madani, and Prof. Mattan Erez were involved in the development of the work presented in Chapter 3.

# TABLE OF CONTENTS

# TABLE OF CONTENTS (Continued)

# LIST OF FIGURES

LIST OF FIGURES (Continued)

## LIST OF TABLES

# Chapter 1: Introduction

---

## 1.1 Near-threhsold error sources and resiliency techniques

The primary goals of this dissertation are to understand the relationship of radiation-induced soft errors and near-threshold operation of digital logic and memory circuits, and to enable and develop new techniques to combat these errors and increase the possibility of reliable near-threshold designs. Prior research has shown that near-threshold operation introduces a significant energy benefit, while also resulting in degraded performance and increased sensitivity to variation. [1]. Little work, however, has been done to understand the relationship between near-threshold operation and soft error rates for both combinational logic and memory.

Circuit designers need to understand the tradeoffs between low-voltage operation, process technology, and soft error rates before chip fabrication. In order to enable them to make these assessments, a simulation framework needs to be developed based on models built using empirical data. To gather this data correctly and consistently, a test platform should be developed that is robust across voltages, and easily portable across processes.

An analysis of this collected data reveals error rate trends that show that low-

ering the supply voltage results in a noticeable soft error rate increase in both the memory and logic circuits. Methods to detect and correct these radiation induced errors as well as errors due to timing variations are determined to be necessary for near-threshold. Additionally, utilizing a massively parallel architecture would help to increase the throughput at low voltages, making the design more attractive for low-power, high-throughput applications. Existing error detection methods, combined with new resiliency techniques to enable efficient error handling on a parallel pipeline are developed, tested, and presented in this work to help advance the concept of implementing practical and reliable low-voltage designs.

## 1.2   Radiation-Induced Soft Errors in Near-threshold

As previously mentioned, a severe lack of experimental research exists that gives insight to the relationship between supply voltage and SER for on-chip memory. Additionally, little experimental neutron and alpha particle works exists for combinational logic at any supply voltage. Knowledge of how the radiation-induced soft error rate increases at low-voltages is essential, especially considering that many applications that require low power also require high reliability. A full experimental understanding of this relationship requires the development of robust and detailed test circuits that can help to characterize SER for many different circuit structures.

In chapter 2 a test platform for radiation characterization across supply voltages is introduced and implemented in a TSMC 65nm CMOS process. This test

chip is intended to be the first in a series of chips across process, that will eventually provide valuable information to circuit designers about the tradeoffs between process, supply voltage, and soft error rate. This chip includes commonly used 6T and 8T memory structures, as well as a variety of logic tests focusing on transient pulse propagation in NAND-based vs. NOR-based logic, the effect of inverter static noise margin on the occurrence of transient pulses, and pulse propagation distance vs. inverter size. Neutron experiments were performed at Lost Alamos Neutron Science Center, and alpha particle experiments were done at the Oregon State University radiation center. The test results presented in this chapter are for experiments that, to the authors knowledge, have not been performed for the particular test circuits across supply voltage.

## 1.3 Synctium-I: A 10-Lane, Near-Threshold SIMD Processor Incorporating Timing Variation Resiliency Techniques

The degraded performance of near-threshold operation, along with the increase in timing errors due to larger variations and increased combinational logic soft error rate makes these low voltage designs the ideal environment for error resiliency techniques. Beyond effective protection against soft errors, these techniques allow for operation at clock frequencies above the timing guard bands, which will result in low-voltage designs with a more acceptable level of performance. Imple-

menting these methods on a parallel architecture will increase throughput even more, enabling them for low-power applications that require a reasonable level of throughput.

In chapter 3 a 10-lane near-threshold SIMD processor is described and implemented in an IBM 45nm SOI process. The processor utilizes two error resiliency methods: lane weaving for static variations, and the Decoupled Parallel SIMD Pipeline to combat dynamic variations. The goal of this design is to increase throughput at low voltages by mitigating the effects of increased variations, while also demonstrating a method that can be used to combat increased soft error rates at near-threshold.

# Chapter 2: Radiation-Induced Soft Errors in Near-threshold

---

## 2.1   Introduction

Operating integrated circuits in the near-threshold regime has become a viable design consideration, due to the significant improvement in energy efficiency [1]. Many emerging applications, such as biomedical devices and unmanned drones, demand low power and therefore could benefit from near-threshold operation. Recent work in this area includes a 64nW ECG SoC for arrhythmia diagnosis [2] and a chip for diagnosis of ventilator-associated pneumonia operating at 0.5V [3]. These devices also require the highest level of reliability, as unexpected errors could halt proper functionality and produce devastating results.

Other applications that are less safety-critical currently experience a tolerable number of soft errors operating as nominal voltages, however as power consumption becomes a higher focus and the operating voltages shift to the near-threshold domain, a certain level of reliability needs to be held to justify the practicality of this change. For example, in the mobile space, recent work includes a 0.48V Pixel-Video recording SoC [4]. Additionally, as the size of data centers increase,

the limitation of power consumption in the data centers while still maintaining high performance is the new design focus. Understanding how the reliability of the systems change is important to understand how ideal near-threshold operation is for these data centers.

Disadvantages of low $V_{DD}$ operation, such as decreased performance and increased sensitivity to variation, are well known. However, little prior work has been performed to quantify and explain the correlation between radiation-induced soft error rate (SER) and low $V_{DD}$ operation for both memory and digital logic. While the basic approach to the relationship between SER and $V_{DD}$ is linear, many simulation models show more complexity based on the circuit response [5]. The non-linear operation of the transistor at low voltage could further complicate its soft error susceptibility and the resulting simulation model.

The goal of this work is to characterize and understand the relationship between circuits operating at a near-threshold supply voltage and their susceptibility to radiation-induced soft errors. Custom test circuits have been developed and tested under neutron and alpha radiation, with memory and logic measurement results that show SER versus $V_{DD}$ that have not been previously performed. This work is the first version of a multi-year project with the intention of fully characterizing the relationship of $V_{DD}$ and SER across processes. While this version will produce unique experimental results for logic and memory error trends across supply voltage, it will also provide valuable information on how to improve upon and expand the test circuits and overall infrastructure of the platform for future generations. The information learned from this test platform, fully portable to

future advanced processes, can help inform future IC designers about the design tradeoffs between circuit design, process technology, low $V_{DD}$ operation, and SER. This information will help them to determine what kinds of circuit techniques will need to be implemented (whether it is more complex error correcting codes or redundancy on memory, radiation hardened memory cells and flip-flops, or error detection and correction on computational logic) in order to ensure the highest level of reliability in low $V_{DD}$ designs.

This remainder of this chapter will be presented as follows. Section 2.2 will cover background information on radiation-induced soft errors by first describing the physical failure mechanism and the concept of critical charge. Results from some initial HSPICE simulations showing supply voltage vs. upset-inducing pulse width will be discussed, followed by a brief introduction of the types of radiation that will be focused on in this work. Existing work in the area of low $V_{DD}$ SER characterization in memory and logic will then be covered. The motivation for this work is given in Section 2.3. Section 2.4 will detail both the memory and logic circuit structures included on the test-chip, giving information about the design and physical implementation of each circuit as well as what conclusions are hoping to be drawn from the data of each test. Section 2.5 will expand the scope to describe the test boards, which are portable to all future generations of the test platform. The top level setup in the radiation environments and a detailed description of the testing procedure will be discussed in Section 2.6. Section 2.7 will cover the experimental results and provide some analysis, highlighting interesting conclusions that can be drawn from the measurements. The chapter will close with

an analysis of the design in Section 2.8, focusing not only on how the conclusions from the measurement results can help circuit designers now, but also on areas of improvement for future versions of the test platform.

## 2.2    Background

### 2.2.1    Charge Collection and Upset Mechanism

Before describing the main sources of radiation affecting integrated circuits, it is important to understand the interaction between the device and the radiation particle that leads to an upset in the circuit. This subsection describes the process of charge generation from an ionized particle passes through the Silicon. The resulting collection of charge at the nearest pn-junction will then be discussed, as well as the state of the transistor that is required to induce an upset in the circuit.

As an ionized particle passes through a material, energy is passed from the particle to that material. The rate (energy transferred per unit length) which this occurs is called the Linear Energy Transfer (LET). This value is dependent on the mass and energy of the particle, as well as the type of material it is passing through. The energy transferred to Silicon manifests itself as charge in the material, where every 3.6eV of energy transferred generates one electron-hole pair in the Silicon [6].

The charge collection process is described in [7], and is as follows. First, a cylindrical track of electron-hole pairs is created as the ionizing particle passes through the material. If the track crosses the depletion region of a pn-junction,

the carriers are separated and charge is collected at the junction via electric-field driven drift and the electrical potential is distorted into a funnel shape. This charge collection mechanism is large and fast, typically occurring within a few picoseconds in modern, submicron semiconductor devices [8]. This results in a large initial current spike, which in many cases is what causes the upset. Once the funnel collapses, or in cases where the particle never crosses the depletion region of the sensitive drain node, charges are either collected at junctions via diffusion, will end up recombining, or diffuse into the Silicon substrate. This will create the long tail on the current profile, or in cases where the particle does not pass through the depletion region it will result in a long, low current pulse that can last up to many nanoseconds.

The amount of charge collected at the junction and whether this can result in an upset depends on a variety of factors. The LET, location, and direction of the incident particle will strongly determine how much charge is transferred to the material, and subsequently collected at the sensitive node. For example, a high-LET particle passing tens of micrometers from the junction will result in a long, low diffusion-only current pulse, where a shorter low-LET particle passing straight through the junction may result in a pulse with a very high initial current spike, and will have a higher likelihood of causing an upset. Whether the type of device hit is a PMOS or an NMOS also matters, as the difference in carrier mobility of electrons and holes will affect how much charge is collected at the junction. At the circuit level, the type of transistor hit also will have an impact depending on whether the switching threshold is above or below 50%. This characteristic will be

discussed in detail later in this chapter.

Other device level characteristics that have an effect of the likely-hood of an upset are whether the the strike happens in the substrate or a well, and if the transistor is 'on' or 'off'. [9] Investigates each of these cases and describes the process of the device response, and if it results in an upset. They find that the outside-the-well 'off' strike is the most sensitive strike location, as it results in the the current profile that is most likely to cause an upset. The drift and diffusion currents raise the struck node voltage and cause an SEU. The outside-the-well 'on' strike reinforces the stored logic state and does not result in an SEU. Inside-the-well 'off' strikes result in a bipolar effect that can cause an upset [10], with smaller gate lengths of more advanced processes only increasing the impact of this phenomenon. For inside-the-well 'on' strikes the bipolar current created tends to restore the node to its original state.

### 2.2.2 Critical Charge

A typical representation of a circuit's susceptibility to an SEU is critical charge ($Q_{CRIT}$), which is defined as the minimum amount of collected charge ($Q_{COLL}$) needed to induce an upset. $Q_{COLL}$ itself is dependent on factors previously mentioned, for example, the particle LET and its direction/location, and the semiconductor characteristics (i.e. doping concentrations and physical geometry). $Q_{CRIT}$ is the most common tool used by circuit designers in assessing a circuits sensitivity to radiation, as increasing the $Q_{CRIT}$ value decreases the probability that the

$Q_{COLL}$ will induce an upset.

$Q_{CRIT}$ is dependent on a number of circuit and device parameters, including which nodes in the circuit are most sensitive (PMOS or NMOS, diffusion area, node capacitance), the ionizing particle types that the circuit is subjected to and its resulting current waveform, and how the particular circuit will respond to a current pulse. To characterize the $Q_{CRIT}$ value of their circuits, designers perform SPICE simulations using a independent current source model developed from 3D device simulations. The simplest model has been proposed by Roche et. al in [11], and is described in Equation 2.1:

$$Q_{CRIT} = C_{NODE} * V_{DD} \tag{2.1}$$

where $C_{NODE}$ is the capacitance of the sensitive node. This approach is insufficient for most cases, as the collection of charge is not instantaneous, and the circuit feedback and response could result in a much different result for a more realistic pulse.

Equation 2.2 shows the Freeman model, developed in [12] and used in [13]. This model incorporates the exponential decay of a time parameter, which is different for each technology and found through device simulations.

$$I(t) = \frac{2}{\sqrt{\pi}} * \frac{Q}{\tau} * \sqrt{\frac{t}{\tau}} * \exp \frac{-t}{\tau} \tag{2.2}$$

In this equation, Q is the total charge deposited by the ionizing particle, and $\tau$ is the process-dependent timing parameter. This model results in a lower magnitude

initial current pulse when compared with the Roche model. However, it also has a much longer decaying tail, providing a more realistic non-instantaneous pulse model that includes the effects of the minority carrier diffusion after the funnel collapses.

The most popularly used model for $Q_{CRIT}$ characterization is the Double Exponential model, which is described in Equation 2.3:

$$I(t) = \frac{Q}{\tau_f - \tau_r} * [\exp \frac{-t}{\tau_f} - \exp \frac{-t}{\tau_r}] \tag{2.3}$$

where $\tau_r$ and $\tau_f$ are rise and fall constants that are process-dependent and determined through device simulations. This equation results in a similar waveform to Eq. 2.2, but with a smaller initial current peak and a wider and longer decaying tail. This current pulse is used for circuit simulations that involve determining the $Q_{CRIT}$ in both SRAM bit cells and single event transients (SET) in combinational logic.

The Diffusion model [14], is intended to model strikes that do not pass straight through the junction, where the pulse is largely due to diffusion current. The model is described in Equation 2.4 as:

$$I(t) = I_{MAX} * [e * \frac{t_{MAX}}{t}]^{\frac{3}{2}} * [\exp \frac{-3 * t_{MAX}}{2 * t}] \tag{2.4}$$

$I_{MAX}$ is the maximum value of the current, and $t_{MAX}$ is the point at which this is reached. As can be expected, this current waveform has a much smaller initial peak than any of the other models, with a much larger and longer decaying tail.

This model – used in combination with the double exponential model – is a good choice for assessment of the $Q_{CRIT}$ sensitivity of low-voltage circuits, as the chance of a long, low pulse inducing an upset may be much greater for the slower devices with lower input thresholds operating in near-threshold.

### 2.2.3   Soft Error vs. VDD Simulations

To better understand how circuits respond to a radiation-induced current pulse as $V_{DD}$ scales, SPICE simulations were performed on 65nm post-layout extracted netlists of a 6T SRAM bitcell (Fig. 2.1a) and a chain of 8 minimum-sized inverters (Fig. 2.2a). The simulation setups use an independent current source to act as a strike on the off NMOS in both test cases. To get an idea of how low-amplitude, long-duration current pulses could affect lower-voltage circuits, a square wave input pulse was used for the simulations, rather than the double-exponential pulse that is commonly used throughout the existing literature. The simulation process involved sweeping the pulse duration with a constant amplitude. The pulse duration was increased, until an upset was observed. This process was repeated for multiple amplitude values at different supply voltages, to see the trend of upset-inducing current magnitude/pulse duration combinations vs. $V_{DD}$.

Fig. 2.1 shows the trend of upset-inducing pulse duration vs. $V_{DD}$ for 5 different current magnitudes for a 6T SRAM bitcell. A medium-length (200ps), low ($5\mu A$) pulse can induce an upset at $V_{DD}$=0.3V, as the pulse-length allows for the output of the first inverter (longer delay due to low voltage operation) to rise high

(a)



(b)

Figure 2.1: Pulse duration vs. VDD simulations for 5 current magnitudes on a 6T SRAM bitcell

enough to cross the switching threshold of the feedback inverter and flip the bit. As $V_{DD}$ increases, the required pulse duration also increases to overcome the increased strength of the cross-coupled inverter feedback. The switching threshold of the inverters also increases to the point where the magnitude of the current pulse is insufficient to induce an upset no matter the length of the pulse duration. The rate

of pulse duration increase as supply voltage increases ranges from incremental (the $40\mu$A magnitude sees an increase of 19ps from $V_{DD}$=0.3V to $V_{DD}$=0.6V) to very large (the $5\mu$A magnitude sees an increase of 85ps from $V_{DD}$=0.3V to 0.45V), with the large increase happening towards the max supply voltage where that particular pulse magnitude can cause an upset. Upsets do not occur at $V_{DD}$=1V until the current pulse reaches $40\mu$A with a 300ps duration, 8x the minimum magnitude necessary to generate an upset at $V_{DD}$=0.3V with a similar pulse length.

The schematic for the logic simulations is shown in Fig. 2.2a. The current source is placed at the output of the first inverter stage (whose input takes a logic zero), while the output of the 8th inverter is observed for identification of SETs. This full setup of 8 inverters actually ended being unnecessary, as it was found that any SET that upset the 2nd inverter would pass through to the end of the chain. The current source reflects a strike on the off NMOS, pulling the node down to zero temporarily. At $V_{DD}$=0.3V and with a pulse magnitude of $5\mu$A, the minimum pulse duration required to generate a transient propagation for a least 8 inverter stages (Fig. 2.2b) is  2x smaller than an upset-inducing pulse in SRAM. As $V_{DD}$ increases, the strength of the inverter driving the sensitive node and a rise in logic switching threshold prevent errors from occurring. While a $5\mu$A pulse can still cause an error at $V_{DD}$=0.3V, a $40\mu$A pulse cannot induce a transient error above VDD=0.75V, showing that combinational logic experiences a significant increase in soft error sensitivity as $V_{DD}$ decreases.

(a)



(b)

Figure 2.2: Pulse duration vs. VDD simulations for 5 current magnitudes on a chain of 8 inverters

## 2.2.4 Sources of radiation

The applications that are the focus for applying what is learned from this project are terrestrial. That is, they are all subjected to the sources of radiation within the Earth's atmosphere. There are three main sources of terrestrial radiation that can affect on soft error rates, all of which are well understood in prior research. These

areas are: high energy fast neutrons (greater than 1MeV), low energy thermal neutrons (much less than 0.001MeV), and alpha particles generated by radioactive isotopes located close to the active areas of the IC's. This subsection will provide background on the sources terrestrial radiation and how they interact with the chip materials.

### 2.2.4.1   Alpha particles

Beginning in the 1970s, alpha particles emitted by trace impurities in packaging materials were found to be a main contributor to DRAM and SRAM soft error rates [15]. The alpha particle is a doubly-ionized Helium atom (consisting of two protons and two neutrons), which is emitted from unstable isotopes such as uranium, thorium, or daughter products in the corresponding decay chains (Po-210, for example). The alpha particle itself is directly ionizing, meaning that it will deposit charge in the devices with an LET of about 0.5 MeV-cm2/mg for a 10MeV particle in Silicon [16].

The alpha particle energies that typically occur range from 2-10 MeV, meaning that the particles themselves only have a range of about $100\mu$m in Silicon. In air, the range of the alpha particles is only 2cm-3cm. Because of this, alpha particles emitted outside of the packaging materials are of little to no concern. Additionally, alphas can be easily shielded, so extra shielding layers on the chip can further limit their SER contribution. Through the use of these shielding layers, along with using purified materials in the manufacturing process, the alpha particle flux rates from

the semiconductor and packaging materials have significantly decreased from 100 cts/cm$^2$/h in the 1970s to about 0.001 cts/cm$^2$/h more recently [7]. The largest remaining contributor to alpha flux is the lead solder on the flip-chip IC, as Pb-210 can decay to Po-210 (a common alpha emitter). This is mostly fixed by moving to lead-free solders, though it has been shown that Uranium and Thorium still exist as alpha particle emitters in the lead-free material [17].

While no longer as large of a contributor to SER as they once were, alpha particles are still important for characterizing the overall reliability of an integrated circuit. For the purposes of this project, alpha tests are used along with neutron experiments for low-voltage radiation characterization. Because an alpha source is also much cheaper and easier to obtain, it is also useful for verifying the test setup functionality over extended periods of test time.

### 2.2.4.2 Neutrons

Unlike the alpha particles, neutrons themselves do not directly create electron-hole pairs in the semiconductor. Instead, they cause soft errors through indirect ionization. That is, the inelastic reactions of the high energy neutrons with the chip materials create a number of different particles that can impart large amounts of charge on the devices. The most common material that is involved in these reactions is Si-28, as it is shown in [18] that accounting for only the n-Si reaction will result in a realistic soft error upset rates (with about a 10%-20% underestimation). As Table 2.1 shows, the products of these inelastic reactions depend on

the energy of the incident neutron. In the past, neutron-induced upset rates were dominated by the burst of energy deposited in a relatively small volume by heavy recoils. However, as $Q_{CRIT}$ values decrease with technology scaling and operation at lower supply voltages, even light nuclear fragments, such as low-energy protons, might deposit sufficient energy in sensitive volumes, contributing to a much sharper increase in SER as supply voltage approaches near-threshold.

Table 2.1: n-Si-28 reaction products and the energy threshold of the incident neutron [19]

| Reaction Products | Threshold (MeV) |
|:---:|:---:|
| $^{25}Mg + \alpha$ | 2.75 |
| $^{28}Al + p$ | 4.00 |
| $^{27}Mg + n + \alpha$ | 12.00 |
| $^{26}Mg + ^3He$ | 12.58 |
| $^{21}Ne + 2\alpha$ | 12.99 |
| $^{27}Mg + 2p$ | 13.90 |
| $^{24}Na + p + \alpha$ | 15.25 |
| $^{15}N + ^{14}N$ | 16.97 |
| $^{12}C + ^{16}O + n$ | 17.35 |
| $^{27}Si + 2n$ | 17.80 |
| $^{26}Mg + p + d$ | 18.27 |
| $^{12}C + \alpha + ^{13}C$ | 19.65 |
| $^{20}Ne + n + 2\alpha$ | 20.00 |

Neutron radiation presents an issue for IC designers largely because, unlike alpha particles, steps (outside of circuit design techniques) cannot be taken to reduce the neutron flux seen by the chip. In order to reduce the flux via shielding concrete needs to be used, where the rate of shielding is only 1.4x lower flux per foot of concrete [20]. Adjusting the manufacturing process will also have little to

no effect, as terrestrial neutrons are created as a result of primary cosmic rays [21] reacting with the top layer of the Earth's atmosphere. The neutron flux is highly dependent on altitude [21], where at 20,000m they reach their peak, and as the altitude lowers, the flux decreases significantly as a result of cascading reactions between the neutrons and the Earth's atmosphere.

The altitude dependency of neutron flux provides a strong motivation for increased reliability of certain applications that operate at higher altitudes. For example, going from sea level to 40,000 ft. (flight altitude for commercial airlines) results in 300x increase in neutron flux. A 5 Mb SRAM that was characterized for 1 error every year at sea level will see about 1 error every 1.2 days of flight time. While the typical passenger is not on a flight for this length of time, passengers with safety-cricital medical devices cannot risk this significant increase error probability compounded with the decrease in reliability from low-voltage operation.

As previously mentioned, both fast neutron (greater than 0.5MeV) and thermal neutron (less than 0.2eV) reactions with the chip materials can result in soft errors in circuits. However, thermal neutrons are well below the threshold energy for Si-28 reactions, and therefore interact with other isotopes found within the semiconductor. It has been found that the Boron-10 isotope found in the Boron-Doped Phosphosilicate Glass (BPSG) dielectric layer was the cause of many reliability issues related to thermal neutrons (Boron-10 has a very high thermal neutron cross-section) [22]. Since this discovery, B-10 has been removed from most advanced process flows, and soft errors due to thermal neutrons have been reduced, though SEU's due to thermal neutrons have still been found in 45nm and 90nm

technologies [23].

## 2.2.5   Related Work

### 2.2.5.1   SRAM

As mentioned in the introduction, a surprisingly limited amount of work has been done to experimentally measure and understand the relationship between $V_{DD}$ and radiation-induced soft errors. Researchers in [24] have done a thorough investigation of neutron-induced SER in SRAM. They find an 18% increase in 90nm cache SRAM SER for every 10% decrease in $V_{DD}$ down to 0.7V, staying well above the threshold voltage. They also look at the dependence of NMOS vs. PMOS diodes and diffusion area. Their results were very interesting, as they found that NMOS diodes had a 14% higher error rate, and also observed a linear relationship between diffusion area and SER.

An in-depth experiment in [25] and [26] of a 10T sub-threshold SRAM in 65nm exposed to both alpha and neutron radiation find an 8x and a 7.8X increase for alpha and neutron-induced SER, respectively, after scaling $V_{DD}$ from 1.0V down to 0.3V. They also take an extended look at Multiple Cell Upsets (MCU) where consecutive bit cells on the same word line are upset by the same particle strike. This is an important effect to take a look at, as an increase in MCUs will require more complex error correcting codes. They find that once the supply voltage drops below 0.6V, the percentage of MCU to total errors begins to increase, topping out

at 3% at $V_{DD}$ = 0.3V. MCUs are highly dependent on bit cell layout and the 10T sub-threshold cell has a much larger layout than the standard 6T cell, and therefore will yield a lower MCU rate at all voltages.

### 2.2.5.2   Combinational Logic

Even less experimental research has been presented for the relationship of $V_{DD}$ and combinational logic SER. In [27] a ring oscillator is implemented in 130nm and monitored for harmonic oscillation, at which an upset event occurs. The threshold energy to induce an upset at each supply voltage is known, as tests are run using the two-photon absorption laser technique [28] and the strike location and energy are controlled. They find that the minimum threshold energy required to cause an upset in the ring oscillator decreases linearly as the power supply voltage decreases, until the circuit enters the sub threshold region at which point the single-event susceptibility remains constant.

While the ring oscillator is an effective method for identifying the threshold energy for SETs occurring on the inverters, trying to determine the error rate by monitoring the frequency changes in an accelerated neutron environment can be difficult. A more effective way to monitor the soft error rate is to have the logic chains output to a storage element, which is then monitored for error counts. This method of combinational logic soft error testing has been done in prior research, however, none of this research has run experiments at low-voltages. These circuits will now be covered, and any necessary changes to perform successful low-voltage

experiments will be discussed.

Researchers in [29] implement the circuit from Fig. 2.3 in 65nm. The test chip contains 4 different chains of 128 D-flip-flops. 3 of the chains have different numbers of inverters in them to see what the combinational logic contribution is at different frequencies ranging from 9.75MHz to 500MHz. They conclude from the heavy ion experimental data that the errors due to upsets in the flip-flops far outnumbered any errors due to logic single event transients being latched. This circuit itself does not focus so much on the soft error rate of combinational logic, as it more just develops an idea of if the combinational logic is as much of a concern as the sequential logic in the design. To shift the focus of the experiment, the sequential elements will need to be hardened to ensure that errors that are being seen are from the combinational logic. Also, the clock buffers will need to be hardened, as a strike on the clock tree could cause a glitch, and subsequently, false errors on multiple bits. This was identified as an issue in [29].

The data path test chip implementation in 32nm for [30] is shown in 2.4. They place 4 different versions of this path on the chip, with inverter chains of length 10 and 6, having both normal P/N ratios and skewed ratios to allow for longer SET propagation. Their on-chip clock generator also can output at speeds from 80MHz to 2GHz. The test chip was was exposed to alpha-particle and neutron radiation, and it was found that the combinational SER increases linearly with clock frequency. It was also found that the combinational SER contribution per sensitive static logic gate is less than 1% of the nominal latch SER at 1GHz at 0.75V. Expanding to the chip level, the total SER contribution of combinational

Figure 2.3: Combinational logic test circuit implemented in [29].

logic is well below 30% of the chip-level nominal latch SER.

Similar to [29], the focus of this work is to identify combinational logic SER trends in relation to clock frequency and to compare the SER of logic when compared with sequential logic. Once again, the FF's were not hardened to isolate the combinational logic itself. In this implementation the clock tree buffers were hardened to prevent unwanted errors from clock tree strikes, which they report to be successful as they saw minimal SER contributions from clock node strikes. While they do not report any supply voltage SER trends, this design as a whole would be acceptable for porting to low-voltages. Increased hardening for the FF's, clock tree, and counters would have to be implemented, possibly by operating them at

Figure 2.4: Combinational logic test circuit implemented in [30].

a separate, higher supply.

## 2.3 Motivation

As the push towards near-threshold design continues, all of the limitations need to be understood and information about them needs to be easily accessible for circuit designers. With more applications emerging where low-voltage operation would be useful, the reliability of these applications and the environment they are in should be well understood. It is important to fully understand the relationship between supply voltage and SER for different circuits within designs.

While some prior experimental research has begun to investigate the dependence on $V_{DD}$ for SRAM SER, their overall assessments are incomplete in different ways. In [24], they do not scale $V_{DD}$ down to the sub/near-threshold region. [25]

and [26] do take measurements down to $V_{DD} = 0.3V$, however, they perform this experimental work on a custom 10T sub threshold cell. The difference in physical cell size of this cell from the standard 6T bit cell may result in largely different MBU rates, which are of great importance. Additionally, the decoupled read and write ability of this cell adds extra capacitance to the internal storage nodes of the cell, resulting in a lower sensitivity to soft errors when compared with the 6T cell.

For digital logic, [29] and [30] present test circuits that are effective for finding the contribution of combinational logic to SER, but do not investigate different circuit characteristics in detail. Furthermore, they do not test at low-voltages, and would need to make some changes to their designs for effective measurements at low supply voltages. [27] uses a ring oscillator to detect SET's, however this focuses on inverter-based logic and this test setup would not be effective for neutron measurements, as they use a two-photon absorption laser with known, controlled energy levels. Most other work found does not involve measurement results, instead performing a simulation-based analysis [31].

The growing interest for near-threshold operation combined with a significant lack of existing research dictates the need for further experimental results in this area. Both commonly used SRAM and combinational logic should be considered. While not ideal for reading to or writing from at low operating voltages, static 6T cells in retention mode at low-voltages will be increasingly susceptible to radiation strikes. Focusing on those, and other more commonly used SRAM bit cells will provide more beneficial and widely-useable results. Sub/near-threshold experimental results for combinational logic will be interesting, and the first of its kind. Specific

results concerning different circuit characteristics also will provide a more in-depth analysis, and be more useful information for circuit designer. Logic measurements across process are also of great importance, as simulations from prior work predict a logic SER increase with process scaling [32].

The goal of the work described in this chapter is twofold. First, measurements taken from the presented test chip are the first of its kind, and will provide useful initial information to the community about some specific circuit dependencies pertaining to the relationship between $V_{DD}$ and SER. Second, and perhaps more importantly, this is the initial version of an entire test platform meant for characterization across processes. The data collected over multiple technologies will then be used to create empirical models and develop a tool for circuit designers to use to determine how SER will factor into their particular design. This first version of the platform will provide useful information for future versions on how the circuits, system, test setup, and test sequence can be improved for more efficient and useful data collection.

## 2.4    Chip Implementation

Because the overall goal of this project is to characterize radiation induced soft errors at low-voltages for both SRAM and combinational logic (particularly taking into consideration different circuit characteristics in the combinational logic), a large variety of circuit structures on-chip are necessary. This raises the issue of test efficiency, as die size and available radiation exposure time can severely limit

the total error count and the resulting statistical confidence of our measured data. Therefore, during chip design heavy considerations were made to maximize the available die space while still developing test circuits that would result in interesting results. Essentially, the goal in the chip design stage was to make the complexity of the test circuits themselves high, while maintaining a simple and robust test-interface.

This section will go into detail about the circuit design choices that were made for this test-chip. It will describe the circuit implementation, discuss why the circuits were chosen, and highlight what information will be drawn from the experimental results. Also noted will be the design decisions that were made to ensure the highest amount of data collection over a limited amount of testing time. Circuit architectures in the test interface to ensure successful low-$V_{DD}$ testing will also be discussed.

## 2.4.1 SRAM

The focus of the SRAM portion of the reference design was to see how commonly used bit cell structures' soft error rates increase as supply voltage is scaled. For these tests, the decision was made to include the standard 6T cell and an 8T cell with a dynamic read port [33]. As mentioned in the previous section, the 6T cell used in commercial off-the-shelf SRAMs is not ideal for low-$V_{DD}$ operation as the low-margins result in read instability and weak write ability [34], [35]. However, the knowledge of the low-$V_{DD}$ SER is still of great use, as many applications will

put the memory cells into a low-power data-retention of 'sleep' mode where the supply voltage is dropped to reduce leakage power.

The 8T bit cell is a popular choice for high-speed or low-power caches, where the read-stability is improved by decoupling the read-line. In [33] they can operate the bit cell down to $V_{DD}$=0.41V in 65nm. Because of this, SER measurements for this bit cell are of use both for applications that utilize a low-voltage retention mode and where the 8T cell is implemented in caches and always operating at low-$V_{DD}$.

A straight comparison to the 6T cells is of use to further assess the merits of increasing memory area by implementing the 8T cells in the design. For example, the effect of added capacitance from the gate of the read-line transistor at sensitive internal feedback nodes should have an effect on per-bit SER. Additionally, different cell layouts (increased size in the 8T cells, different location of sensitive nodes) also should have an effect on multiple bit upset (MBU) measurements.

The high-level implementation of the memory test structures is shown in Fig. 2.5. The structure of both the 6T and 8T bit cell arrays are shown on the right portion of the figure. Portions of the circuit drawn in black are included for both the 6T and 8T structures, circuits drawn in red are included only for the 8T, and dotted blue lines are included only on the 6T section. For both structures, each column contains 32 cells on a shared bit line, with well taps placed on both the top and bottom of each column. The 6T and 8T arrays have 560 and 462 cells, respectively, on a shared horizontal word line. The left portion of Fig. 2.5 shows the block diagram of the memory floor plan. Two banks of memory arrays share

Figure 2.5: Top-level SRAM architecture

one synthesized scan-chain block, such that each 6T and 8T scan chain block is 1120 and 924 bits, respectively. This grouping was copied 3 times for both the 6T and 8T cells to maximize the number of bits on-chip limited by the available die space.

Writing to/reading from the bit cells is done through a dual-clocked scan in-

terface. The dual-clocked scan prevents any race conditions from occurring, where the data could pass through consecutive scan-bits on one rising clock edge, resulting in lost bits. The scan chain schematic for each column of memory is shown on the bottom right portion of Fig. 2.5. The input to the first flip-flop (clocked by clk) takes either the read line value (the non-inverted bit line for the 6T array or the decoupled read line for the 8T array) in parallel mode, or the output of the previous bit in the scan chain in serial mode. This is determined by a multiplexor whose select input is controlled by a 'scan_enable' signal that comes from off-chip. The output of the first flip-flop for each bit of the scan chain is connected to one column's bit lines, with the inverted value connected to the complemented bit line of each column in the memory array. Two tri-state buffers controlled by an enable bit (WEN) prevent the scan-chain input from incorrectly affecting the state of either bit line.

As shown on the left side of Fig. 2.5, each memory array bank has its own local hardened address decoder. All address decoders take the same address select inputs from an off-chip source. For word line activation, the desired address bits are input directly into the address decoder, which is hardened through logical masking by simply ANDing the decoder output with an enable signal to prevent an address from being incorrectly accessed during testing. This setup, while simple, is efficient for the types of test that we want to run, as a known set of bits will be clocked through the scan-chain and written to all of the columns simultaneously for scan-in, and each word set will be loaded into the scan-chain and then clocked through on scan-out. Both of these procedures will only need to be repeated 32 times, once

for each bit in the column.



(a)



(b)

Figure 2.6: Memory bit cell custom layouts: (a) 6T; (b) 8T

Due to legal issues with placing the TSMC high density 6T cell provided by the foundry into our design, both the 6T and 8T SRAM bit cells utilized custom layout based on logic design rules. This resulted in physical cells (Fig. 2.6) that

were roughly twice as large. The large spacing had three major effects on the tests. The first, and most detrimental, was that the amount of memory that could be included on chip was cut in half, significantly affecting the test efficiency. Second, the custom cells were designed with the cross-coupled inverters having equal P/N ratios, meaning that the static noise margin of the inverter is shifted from %50 and will have an effect on the error rate (as a lower switching threshold should be more sensitive to current pulses). The third issue was that the lower density of the memory arrays – with sensitive internal nodes that were spaced farther apart – would lead to MBU results that were much lower than they would be for the standard 6T bit cell arrays. Using larger custom designed cells had a few positives, however, as this allowed for more control of internal node capacitance and the ability to monitor the physical layout effects on MBU's. For example, it can be seen in Fig. 2.6 that the PMOS and NMOS are isolated from each other to the left and right. Monitoring MBU trends for adjacent horizontal cell orientation will allow for observation of the effect of the placement of sensitive PMOS or NMOS drain nodes next to each other in local bit cells.

## 2.4.2   Digital Logic

The focus of the digital logic tests was to identify how particular implementations of combinational logic are affected by radiation-induced soft errors as the supply voltage is scaled. In order to achieve this, a variety of different test circuits on chip are necessary, all needing to be fully functional at low voltages. This requires

steps to be taken to increase test efficiency and accuracy. This subsection will describe the methods used to implement the combinational logic experiments. It will first describe the high-level architecture and test interface, and then go into more low-level circuit details about each of the tests.

The high-level architectural diagram of the logic test structure is shown in Fig. 2.7. Each individual logic test consists of the following: The circuit under test, a level shifter, and a 5-bit counter. There are 20 unique logic chain implementations in each set – all with a synchronous and an asynchronous version – for a total of 40 different chains. Asynchronous tests were used for strict observation of the tests themselves, that is, any SET that occurs in the chain (As long as the pulse width is long enough) will be counted, which shows the direct effect of each circuit characteristic on whether or not SETs are occurring in the first place. The synchronous tests were placed with the purpose of testing the effect of timing masking and how the error rate changes with clock frequency. The timing masking effect is important when considering increasing error rates as voltage is scaled, as a lower clock frequency in the sub/near-threshold region will prevent many SETs from being clocked into flip-flops and actually manifesting themselves as errors.



Figure 2.7: Top-level digital logic architecture

Every set of 40 tests is duplicated 10 times (the maximum allowed by available die space) to increase the probability of an SET occurrence during testing and achieve a large enough sample size of errors. Even with this increase in tests, each test only totaled, at most, 320 logic gates on-chip. Recent work has shown that a standard sized inverter in 90nm has a FIT (Failures in Time – the number of failures for every billion hours of device operation) rate of $4.4 \times 10^{-4}$ [36]. Comparing this values with a standard memory 90nm FIT rate of $1 \times 10^{-3}$ gives a static logic SER per gate about 45% lower than the SER per memory bit [37]. Considering that there is 300x fewer cells per logic test on the chip when compared with the 6T array, it will take over 600x the amount of test time to obtain similar error numbers. While this was not ideal for the neutron tests and will need to be fixed for future test chip generations in this project, it was found that enough time was provided with the alpha tests to observe a sufficient number of errors.

A 5-bit counter is used for each chain to detect a maximum of 32 errors during each test. The very short length of each test sequence essentially negated the chance that any more than a few errors could occur during any single run. Each of these counters is radiation-hardened through two methods. The first is to operate the counters at nominal voltage at all times no matter what voltage the circuit under test is operating at, and the second method is to use custom built DICE flip-flops (Fig. 2.8) for each bit. The DICE [38] storage element is a unique cell that is hardened by design, where a strike on any one node of the keeper portion of the cell will not change the value stored in the cell. A strike that upsets two nodes could still induce an upset, however, so steps were taken during the custom

physical layout to minimize the chance of a single particle strike upsetting multiple nodes. Additionally, any errors occurring on the counters were detectable through monitoring the change in error count over the sequence of scan-outs. For example if the data was scanned-out every 10 seconds, an increase in the counter greater than 2 was considered a result of an upset occurring on a flip-flop in the counter. Keeping with this methodology would result in only minor differences in the overall error counts.



Figure 2.8: Radiation hardened Dual-Interlocked Storage Cell (DICE) [38] used in the 5-bit counters

Each bit in the counter has a partner bit in the scan chain (which consists of unhardened flip-flops operating at a nominal supply voltage), so that the length of the logic scan chain is the total number of logic tests times 5-bits. Similar to the memory tests, TMR was used on the scan chains to detect any errors that may have occurred during scan out and each redundant scan-chain has its own input and output pin, to reduce test scan-in/scan-out time. Identification of scan-

chain errors was done during post-processing using a majority voting method. To prevent the possibility of an inadvertent particle strike on the reset line, the reset ability for the latches in the counters was removed. This eliminated the capability of zeroing out the counters prior to each test period. To account for this during the experiments, it was necessary to scan out from the counters twice: once before irradiation, and then again afterwards. Subtracting the initial value from the final value gave the total number of errors that occurred during that test period.

To ensure that an SET in the logic chain will always latch into the counter, no matter what voltage the logic under test is operating at, a buffer – serving as a level shifter – is placed between the end of the logic chain and the counter. This buffer is set to its own supply voltage ($V_{LS}$) halfway between the logic chain voltage and counter voltage, and sized 2x minimum to properly drive the input of the DICE flip-flop. In order to account for the possible contribution of the level-shifter buffer on SETs, duplicate tests consisting of only a level-shifter connected to a counter are included. Any errors that occurred on these chains are then subtracted from the total of all other logic tests. No errors were actually detected for any of these level shifter circuits, most likely because the increased size of the cells used in the level shifters when compared with the logic under test.

As previously stated, all logic tests are duplicated to include both synchronous and asynchronous versions. The lone difference between the two versions is that the synchronous test places a clocked D-FF between the level shifter and counter in each logic chain, so that only errors caught on the rising clock edge are counted. The clock tree in the synchronous tests is hardened by increasing the size of the

buffers by 32x. The counters in both tests consist of toggle flip-flops, where the input comes straight from the level shifter in the asynchronous test, and from the clocked flip-flop in the synchronous test.

The following subsections will detail the individual circuit implementations of the digital logic tests. A variety of unique logic chains were implemented on-chip to isolate and analyze the VDD vs. SER relationship of different circuit characteristics. These characteristics include: inverter static noise margin (SNM); NAND-based vs. NOR-based logic (focusing on differences in the cell design that can lead to higher SER); and inverter size vs. transient pulse propagation. An important aspect of all of these tests is that they use a fully digital implementation with as many cells from the TSMC standard cell library as possible, thus making these tests easy to implement for future test chip versions.

### 2.4.2.1   Static Noise Margin

The purpose of the first test is to observe the effect of the switching threshold of an inverter on the probability of an SET occurrence. The conclusions drawn from this experiment are twofold. The first being an observation of the trend of how the error rate increases vs. increasing noise margin, a determination that can inform the circuit designer of how they can alter their logic to take advantage of the soft error protection provided by increasing the switching threshold. The second is to gain an idea of how effective increasing the noise margin is as a technique to use when the supply voltage is scaled. That is, similar radiation pulses will occur at

any supply voltage, so a pulse that induces an upset in an inverter operating at $V_{DD}$=1.0V with a 50% SNM should also induce an upset in an inverter operating at $V_{DD}$=0.5V, no matter what the SNM is. If the noise margin needs to be increased far above 50% under low-voltage operation this technique could become impractical, as the propagation delay of the inverter could increase significantly.



Figure 2.9: Static noise margin: (a) test chains; (b) simulated switching thresholds

Fig. 2.9a shows the schematic for the SNM test implementation. Each chain

consists of 32 inverters symmetrically sized to maintain identical switching thresholds within each node of the chain. All inverters in this test utilized a custom layout, for full control of the noise margins. The cells were then placed in the digital flow and inserted into the gate-level net list. The margins tested ranged from 20%-80%, in increments of 15%. A static input supplied by an off chip source sets the switching threshold (for 1 to 0 or 0 to 1) for each test. The simulated input thresholds with a static input value swept from 0 to 1 for each test are shown in 2.9b. A static input value of 1 reverses the margins of the tests so that the observed error rates should be the same for opposite tests, i.e. the error rate for 'SNM 1/1' with 'static in' = 0, should be the same as 'SNM 8/1' with 'static in'=1. For both cases, the error rate for 'SNM 2/1' should remain the same.

## 2.4.2.2   NAND-based vs. NOR-based Logic

This second combinational logic test looks at what approach to take when choosing what base cells to use for implementing digital logic circuits. NAND and NOR gates both exhibit the property of logical completeness, which means that any boolean function can be completed using either all NAND or all NOR gates. Because of this, often times in digital circuits are constructed strictly of NAND or NOR gates. Typical implementation at high speeds is done using NAND gates, as the propagation delay of the NOR-gate is higher. However at lower supply voltages where the operating frequency decreases significantly, using NOR-based logic is a more reasonable idea given its lower area and energy when compared with NAND

gates.

The details of the NAND/NOR tests are shown in Fig. 2.10. Each test is 32 gates, with the same input of every gate in the chain held to the same value to emulate the functionality of an inverter chain. There were three unique tests in this set included on the test chip. The first is 32 NAND-gates with the 'B' input held to a static '1', while the 'A' input takes the output of the previous gate, with the first gate in the chain taking its input from an off-chip source. The second test (labeled NOR-X in Fig. 2.10) is 32 NOR-gates in the same orientation as the NAND test, but all with their 'B' inputs held to a static '0'. The NOR-Y test switches the inputs from the NOR-X test, so that all of the 'A' inputs take a static '0', and the 'B' inputs take the output of the previous gate.

The bottom portion of Fig. 2.10 shows the transistor level schematics of the gates used in the test. Transistor sizes are indicated as a P/N multiplier value. Nodes that are most sensitive to a particle strike (the drain of an off transistor) at any point within the chain are identified and highlighted in red.



Figure 2.10: NAND and NOR test chains and radiation sensitive nodes.

Comparison of the NAND and NOR-X chains focuses on the differences in node capacitance and diffusion area of the sensitive nodes. At the only sensitive node in either chain, the NAND gate has two 6x PMOS transistors and one 5x NMOS transistor, whereas the NOR-X gate has one 6x PMOS transistor and two 4x NMOS transistors. So while the NAND gate has a larger total capacitance, it also will have larger diffusion area, implying that the results of this test will reveal which characteristic has a greater effect on SER, as there will always be a tradeoff of diffusion area versus node capacitance. Assessing the NOR-X and NOR-Y tests looks at the different locations of sensitive drain nodes throughout the chain. The NOR-X test only sees the output node of each gate as the sensitive node, where the NOR-Y will see two separate sensitive nodes when the 'B' input is 1 and the 'A' input is 0. The effect of this is twofold, as not only will having two sensitive nodes in one gate increase the probability of upset, but the lower node capacitance seen by the drain node between the stacked PMOS transistors will also increase the error probability of the NOR-Y configuration.

### 2.4.2.3   Inverter Size/Chain Length

The final logic test observes the effect of inverter size on pulse propagation as supply voltage is scaled. The schematic for this test is shown in Fig. 2.11. Separate chains of length 2, 4, 8, 16, 24, and 32 were implemented for inverter sizes of 1X, 4X, and 8X. Pulse propagation length was determined for each inverter size by monitoring the error rate vs. chain length trends. If the error rate continued to increase at

43

a significant rate as the chain length increased, it was assumed that errors were continuing to fully propagate through the inverter chain. Once the error count plateaued, the conclusion was made that SETs could not propagate fully through the chain. By implementing the full set of varying chain lengths for 3 different inverter sizes the trend of gate size versus SET propagation could be determined, and an idea could be developed of how large inverter chains need to be sized to restrict SET propagation to only a few logic gates at low voltages.

Figure 2.11: Inverter chain size/length test schematic

### 2.4.3    Physical Implementation

The full layout for the initial version of the SER reference test chip is shown in Fig. 2.12. The chip has 72 total pins, with 9 different power domains: logic, counters, level shifters, SRAM, pulse measurement diode, and 4 different scan chains. There are 12 total scan chains on the chip (6T, 8T, 10T, and digital logic). The SRAM and pulse measurement portions of the chip utilized full custom layout, while the digital logic portion was a fully automated digital flow. Synthesis for this digital portion was done using the Synopsys Design Compiler, and place and route was done using Cadence Encounter. All custom layout, as well as the top level chip implementation was done using Cadence Virtuoso, where final DRC (Design Rules Checks) and LVS (Layout Versus Schematic) checks were performed.

Up to this point in this work, the portions of the chip labeled 10T SRAM and charge collection and pulse measurement on the layout in Fig. 2.12 have not been discussed. The 10T SRAM is an interruptible bit cell that has much improved write ability and read stability, which makes it an appealing choice for low-voltage designs. This cell was included in the SRAM tests, however, the bit cells were not designed properly, leading the unreliable test results. Therefore, the cells were not included in the design details, and these results will be omitted from the analysis portion of this thesis.

The charge collection and pulse measurement test is a variety of different types and sizes of diodes whose amplified outputs are sent off-chip to an oscilloscope. The goal of these tests was to capture the waveform of current pulses resulting

Figure 2.12: NAND and NOR test chains and radiation sensitive nodes.

from ionizing particle strikes on the diodes. Capturing these waveforms would provide a direct measurement to verify the current pulses that are used for circuit simulations in SPICE. Our tests were unsuccessful, as we were unable to capture any waveforms that we could conclude without a doubt were a result of ionizing radiation. Further details about this test, as well as what should be changed to make this experiment successful in future versions will be discussed in the design analysis section of this chapter.

## 2.5 System Implementation

This section will describe the design of the test-system used during both the alpha source and accelerated neutron beam experiments. The board-level implementation involves the combination of one master board and as many as 10 sub boards. These boards will both be discussed in this section, highlighting any unique techniques used to increase the test efficiency and system robustness.

## 2.5.1 Master Board

The master board provides the main interface between the sub boards (which contain the actual devices under test) and the test control FPGA hardware. The board (shown in Fig. 2.13a is 15 inches by 6 inches and has 10 PCIe slots, which the sub boards plug into vertically. The PCIe is beneficial in that it provides plenty of signal pins to properly operate 8 chips up to 10GHz on each sub board. It also gives the sub boards a vertical orientation that is suitable for beam operation without needing the assistance of any external equipment.

Lower speed signals (including signals for the scan interface) are delivered to the board from the FPGA using a 40-pin VHDC connector. High-speed clock signals are generated using an on-board PLL chip [39], which sends a separate clock output to each sub-board. All connections for the PLL (including control signals and power/ground signals) come from the FPGA board, through low-speed headers. The 9 on-chip voltage domains are consolidated into 4 groups that take the same supply voltage (scan-chains, devices under test, level shifters, minus input

for diode tests), with each group's input provided by a separate power supply and shared between all of the sub-boards.



(a)



(b)

Figure 2.13: Master board: (a) image; (b) block diagram, showing the scan signal routing jumpers.

The most unique design aspect of the master board is the large array of jumpers duplicated for each sub board interface. This setup is shown in Fig. 2.13b. Each low-speed scan input/output signal is routed off of each sub board, through a set

of jumpers on the master board, and back into the sub board. These jumpers allow the user to route any scan-chain input and/or output signal around any individual chip. This ensures that if any part of any chip ends up with a functional failure, only that section of that chip has to be removed from testing. This maximizes the data still collected from the other tests that remain operational on that chip. All other chips remain unaffected by this process, and proceed to function as normal.

### 2.5.2 Sub Board

Each sub board (Fig. 2.14a) in the test setup has 8 chips oriented symmetrically equidistant in a 2 inch diameter circle at the center of the board, to be aligned with the 2-inch collimated neutron beam. Each chip was placed and wire bonded directly to the board, since the use of any packaging would have prevented all of the chips from fitting. Future versions of this platform should use a re-design of this portion, where there are fewer chips in the center (either by placing some chips on the back of the board, or omitting some of them altogether). This will allow for the use of packaging (much cheaper and less difficult to implement than chip-on-board) and decrease the routing complexity in the center of the sub board, which ended up being the bottleneck for the board design and cost.

The block diagram of the sub-board implementation is shown in 2.14b. Each chip on the board has a set of jumpers that allow the user to disconnect each individual supply voltage from each chip separately. As with the scan-signal jumpers on the master board, this is useful in the event of a functional failure of part of

(a)



(b)

Figure 2.14: Sub board: (a) image; (b) block diagram.

a single chip, as it allows for the removal of any unnecessary switching noise or shorting on the voltage lines (due to wire bonding errors) from unused chips. Each chip also has its own scan input/output through the PCIe interface to allow for the scan routing to be done at the master board level. High-speed clock distribution is dealt with locally using a clock buffer on each chip [40]. The clock buffer takes the reference input from the master board PLL, and has up to 8 differential outputs used as single ended clocks. Additionally, there are 8 separate SMA outputs for the pulse capture diode tests. Each chip outputs to one SMA (with each chip testing a different size/type of diode), which connects directly to a high-speed oscilloscope. While the on-chip portion of the diode tests themselves were unsuccessful, the board-level design for delivering the signal to the oscilloscope was found to be effective.

## 2.6   Top Level Setup

This section describes the test setup decisions made at the highest implementation levels. The section will begin with details about the equipment used for the top level setup and organization of the equipment and boards for the accelerated neutron and alpha tests. Then the FPGA implementation, including discussions about the software and data interfaces, will be covered briefly. Lastly, a description of the test sequence will be given.

## 2.6.1 In-beam Test Setup

Neutron tests were executed at the Weapons Neutron Research Facility (WNR) at the Los Alamos Neutron Science Center (LANSCE), in the Irradiation of Chips Electronics (ICE House) beam room [41]. The scan input/output, control signals, and PLL control were all dealt with using the Digilent Genesys development board [42]. This board contained a Xilinx Virtex V FPGA, running a microblaze IP processor [43]. The FPGAs processor was then programmed in C.

The test setup in the beam room is shown in Fig. 2.15. All tests were initiated using Python code from a control computer on the safe side of the 14 foot high concrete beam room wall. The computer then communicated, through Ethernet, to the Genesys development board and a GPIB to Ethernet switch, which ultimately controlled the power supplies through GPIB. Four outputs from the power supplies where then connected to the main board using standard banana cables. All of the 1V I/O and control signals to/from the chips were sent using a VHDC connector, while the 3.3V signals for controlling the PLL were sent through SPI. Even though the power supplies operated from inside of the beam room, they were well away from the beam to the point where neutron scatter from the beam did not affect the functionality of the parts. The FPGA board also was located inside the beam room, just 6 inches from the beam. This was found to be common practice, and throughout the testing no errors were observed on the FPGA.

Fig. 2.16 shows an alternate view of the test setup in the beam-room. The setup was mostly standalone, however, it needed to be placed on risers to become

(a)



(b)

Figure 2.15: Beam room test setup: (a) block diagram; (b) image

Figure 2.16: Neutron beam test front view

level with the beam. While the system was designed to accommodate up to 10 boards in the beam, only 3 were tested at a time due to functionality, and power integrity issues. The beam width was 2 inches, and the sub boards were aligned with the beam using a laser sight. The test bench was shared with 2 other groups.

Tests were run at Los Alamos Neutron Science Center (LANSCE) for 9 days, 24 hours per day, during August/September, 2013. Error rates were calculated by monitoring the number of counts on a pulse counter located at the beam source. The rate of neutrons per counting pulse changed daily and the translation factor was given to us the following day by the LANSCE staff. Over the 9 day period, the neutron flux ranged from 5e5 to 9e5. These rates resulted in an average of about 4 hours of testing time for each voltage and input sequence, to accomplish

our goal of 1000 errors/Mbit of SRAM.

## 2.6.2   Alpha test setup

Alpha experiments were run at the Oregon State University radiation center for 3 months. The tests utilized an Am-241 point source, which has a 2-pi surface emission rate of 1.16e5 alphas/minute, and an activity of 0.1098 uCi measured in 2000. For the most part, the same overall test setup as the neutron experiment was used. However, there was one main difference, which had an effect on data collection efficiency. This difference was that the alpha source was too small to cover more than 1 chip, and according to JEDEC test standard JESD89a [44] the source must be placed as close to the bare face of the chip as possible. To properly setup this experiment, only 1 sub board was used with a petrie dish placed upside down over the chips. A small hole was cut so that the Am-241 source fit tightly, and was only 2 mm above the chip (Fig. 2.17). Because only 1 chip at a time was tested, data collection was much slower. However, having 3 months to run tests allowed for a large enough amount of collected data to result in a high statistical confidence.

## 2.6.3   FPGA software architecture

The hardware and software running on the FPGA has been specifically designed to transfer the large amounts of data that are being sent to and from the test chips

Figure 2.17: Alpha experiment setup

easily with a host computer. The FPGA runs a light-weight webserver allowing for simple SET and GET commands to be performed from any web-enabled computer. Once a specific URL is visited on the FPGAs webserver, a respective C subroutine is called executing a predefined action such as scan in, scan out, enable PLL, etc. All commands return a true/false token based on the success of the command and can easily be parsed in Python by the host computer.

In the case of scanning data in, both SRAMs on all chips can be scanned in simultaneously, enabling an increase in efficiency over a system where paralleliza-tion is not an option. Because SRAM input data is generated by the FPGA (i.e. all-1s or all-0s) it does not need to be stored before transmission. This type of action also requires no data to be returned.

An operation such as scanning out does require stringent data-storage require-

ments. Scanning out large amounts of data cannot simply be streamed through a TCP connection because the SRAM on the FPGA cannot buffer the data as the FPGA scans out from GPIO faster than Ethernet can transmit. This is solved by saving the scanned out data directly to an onboard file system in DRAM. This also enables faster testing time because the scanned-out data be retrieved from DRAM at a later time after all scanning (out and back in) has completed. Data is later retrieved by simply pointing to a specific URL and downloading a file through a reliable TCP socket.

### 2.6.4   Test sequence

The flow of the test sequence used during neutron tests is shown in Fig. 2.18. $V_{HARD}$ is the supply voltage used for all portions of the design that require as much radiation hardening as possible. This includes the high I/O voltage, the logic counters, and the SRAM address decoder. $V_{DUT}$ represents the supply voltage for all of the circuits that are being tested. This includes the SRAM bit cells, logic test chains, and scan-chain flip-flops.

   The test flow is shown in the color coded boxes of Fig. 2.18, and is as follows: Prior to beam exposure the logic counters are scanned out. Then, a pre-determined pattern (either all 1s, all 0s, or checkerboard) is scanned into the SRAM arrays and scan flip-flops. The voltage of the circuits under test (SRAM arrays and digital logic tests) is then lowered. The beam is then turned on, and the logic tests are scanned out over short intervals, to ensure that much less than 32 errors per test

Figure 2.18: Neutron beam test sequence

over a single run are collected. After the end of the beam exposure period the number of pulse counts are recorded for error rate calculation purposes, the test circuit voltage is raised back to 1V, and all of the tests are scanned out.

Baseline scan-out tests were run with the beam on and it was found that no errors occurred during the duration of the scan-out for many attempts. While the SRAM and flip-flop tests were static, and did not require any scanning during beam exposure, this baseline determination was important for the logic tests, as it was necessary to scan-out about once every minute. Despite this, any errors that could have possibly occurred on scan-out during active tests were detected and corrected using the scan chain TMR. Since the baseline tests deemed the TMR unnecessary, removing the TMR for future versions is a reasonable choice that would allow for much more die space for the circuits under test.

## 2.7    Measurement Results



Figure 2.19: Die Photograph

Fig. 2.19 shows the die photo and Table 2.2 is design summary of the 2mmx2mm test chip, fabricated in a 65nm CMOS process. All tests were operational for supply voltages ranging from 0.33V-1.0V. The chip contains 107.5 kb of standard 6T cells and 88.7 kb of 8T RF cells, both with a custom layout. There were also 105kb of 10T cells, though those were omitted from inclusion in the table since there are no results reported for these cells. While there were 40 total logic tests, the high speed clock generation from the off-chip PLL was not debugged completely prior to radiation testing, and no measurements were taken for the synchronous tests. Therefore, results are reported for 26 logic tests (all asynchronous), each dupli-

cated 10 times, for a total of 260 test logic chains. It should also be noted that memory results are reported both for neutron experiments at LANSCE as well as alpha experiments at OSU, while the digital logic test results are reported only for the alpha particle experiments due to a cross-section for the logic tests that was not large enough to result in enough errors in the neutron beam for statistical validity.

| Process Node | 65nm CMOS | |
|---|---|---|
| Die Size | 2mm x 2mm | |
| VDD | 0.33-1.0V | |
| # of Transistors | ~3.3 million | |
| Memory | | |
| Cell | Cell Size | Array Size |
| 6T | 0.849μm$^2$ | 107520 bits |
| 8T | 1.173μm$^2$ | 88704 bits |
| Digital Logic | | |
| # of Tests | 26*10=260 | |
| Block Size | 750μm x 1160μm | |

Table 2.2: Design summary of the soft error test chip.

## 2.7.1 SRAM

Fig. 2.20 shows the trend of SER vs. retention $V_{DD}$ for the 6T and 8T memory arrays for both neutron and alpha radiation. Determination of the minimum voltage was made by running baseline tests outside of the radiation environment, where all of the cells were first written-to, and then read-from. $V_{DD}$=0.33V was found to be the minimum supply voltage required for reliable data retention. From

$V_{DD}$=1V down to $V_{DD}$=0.33V, the 6T SRAM SER increases by 6.45x for acceler-

ated neutron radiation (Fig. 2.20a). For a supply voltage decrease from $V_{DD}$=1V

down to $V_{DD}$=0.35V, the 6T SER increases by 2.5x for accelerated alpha radia-

tion (Fig. 2.20b). Two different data sets are shown in each plot. The 0-to-1 data

curve corresponds to the case where all 0's are scanned into the SRAM prior to

irradiation. Cells that read a 1 after scan-out are counted as errors. The 1-to-0

curve represents the opposite case, where all 1's are scanned in and any 0's after

scan-out are counted as errors.



Figure 2.20: SRAM/RF SER vs. Retention $V_{DD}$ for (a) neutron and (b) alpha
radiation.

The reason for the difference between neutron and alpha trends lies mainly

in the nature of the two types of radiation, as well as the overall distribution of energies. The alpha particles are emitted mono-energetically (all with the same LET) from the Am-241 source, and mostly physically interact with the sensitive nodes in the same way (passing through the junctions in a similar direction). This results in a limited distribution of collected charge at the junctions, such that while lowering the supply voltage increases the cell's sensitivity to lower amounts of collected charge, there is not a significant amount of increasing instances of ionizing particles passing through/near the junctions and accumulating the lower amounts of charge. Neutrons, on the other hand, have much more uncertainty for a number of reasons. The spectrum of incident neutron energy is much larger (1MeV-800MeV in the ICE House I beam at LANSCE), and the reactions of these neutrons over this spectrum of energies with the Silicon results in different products (Table 2.1) that have unique LET values. Additionally, the products of these reactions emit at different angles, resulting in the ionizing particles taking different paths through/near the sensitive junctions. For these reasons, the distribution of collected charge is much larger, resulting in a larger increase in upset probability at lower supply voltages.

Once $V_{DD}$ drops below the threshold voltage of the transistors ($V_{TH}$ = 0.35V) the accelerated neutron error rate increases sharply. This is likely a result of the transistors entering the weak-inversion region, where $I_{ON}$ has an exponential response to any change in gate voltage, increasing the cell sensitivity significantly. Also, as was shown in Fig. 2.1, lowering $V_{DD}$ down to the threshold voltage exposes it to small magnitude, short duration pulses. Further simulations show

that decreasing the supply voltage below the threshold voltage, even just to 0.33V, increases the sensitivity to an order of magnitude smaller current pulses.

The 8T error rate remained roughly similar to the 6T trends for all cases (neutron or alpha radiation, 1 to 0, or 0 to 1), and at most test points, the error rate was even slightly smaller than the 6T data. This is to be expected, as the added gate capacitance from the access transistor on the dedicated read line should slightly lower the SER of the 8T cell. This difference should be small, as the rest of the cell (most importantly, the cross-coupled inverters) is identical between the 6T and 8T. One case (neutron data, 1 to 0) showed higher SER for the 8T cells, which goes against what is expected. However, looking at the data shows that the difference is within the statistical error bars, as a lower sample size from the neutron test may have affected the final measured error rates.

The neutron-induced multi-bit upset (MBU) rate vs. $V_{DD}$ for the 6T SRAM is shown in Fig. 2.21. The plot shows the contribution of three different types of upsets (single bit (SBU), double bit (DBU), and triple bit (TBU)) to the total number of errors, in terms of percentage. For reference, there were 2000 total measured errors at each voltage (1000 errors/Mbit of SRAM). MBUs are considered adjacent cells on the same wordline that are upset by the same ionizing particle strike, as inter-wordline errors require increased complexity in the error correcting codes (ECC) (and the existence of TBUs increases the ECC complexity even more). While the likelihood exists that two adjacent cells were upset by two separate ionizing particles during irradiation, the probability was low – as 1 hour of testing averaged 500 errors for 2 Mbit of data – so these MBU occurrences were always

assumed to be from the same ionizing particle strike.



Figure 2.21: SRAM MBU rate vs. $V_{DD}$ under neutron irradiation.

Fig. 2.21 shows that DBUs increased from 2.5% at $V_{DD}$=1V to 6.4% at $V_{DD}$=0.33V, while TBUs increased from 0% to 0.2% over the same supply voltage range. In total, MBUs experience a 2.6x increase when $V_{DD}$ scales from 1.0V to 0.33V. No TBUs occurred until $V_{DD}$ reached 0.5V, increasing after that point by 2x when lowering the supply voltage down to 0.33V. Note that, as previously mentioned, the custom SRAM cell size obeys logic design rules, such that the area is 2x larger than a commercial memory. This will have a strong effect on TBUs, as one ionizing particle will have to pass through three cells horizontally, and having 2x larger distance to pass through will significantly decrease the probability of all three cells being upset.

Taking a deeper look into multi-cell upsets (MCU) – adjacent cells (wordline or bitline) that are upset by the same particle strike – reveals that they are heavily

dependent on the physical layout. 98% of MCUs were adjacent bits on a horizontal word-line with sensitive nodes 0.37$\mu$m apart. The other 2% (shared bit-line) were separated by 0.8$\mu$m.

Additionally, there was a large difference in MCUs that were observed depending on column orientation, and which transistor type appeared to be most sensitive. 90% of shared word-line MCUs occurred over a mirrored column orientation where the sensitive areas of adjacent bits were NMOS drain nodes. It is possible that these NMOS nodes are more sensitive due to their inherent device properties. Electron mobility (the majority carriers in n-type material) is greater than hole mobility (majority carriers in p-type material), therefore the diffusion length of the majority carriers in the NMOS diffusion is longer, thus having a greater chance to collect at the junction. Also, since it is a p-type substrate, the PMOS diffusion area sits inside of an n-well, so any carriers created by a particle passing through a shared n-well will migrate towards the well-ties rather than collecting at the sensitive junctions. Carriers resulting from ionizing particle passing adjacent NMOS diffusions in the p-type substrate have a better chance of reaching the junctions, as they will be farther from the substrate taps.

Approaching a conclusion from strictly a device physics standpoint is not entirely appropriate, however, as the circuit response also has a lot to do with an upset occurring. To save as much area as possible, the P/N ratio in the cross-coupled inverters was set as 1/1. This moved the switching threshold away from 50% towards where 1 to 0 transitions happen at a lower input voltage. The result of this is that, when the drain of the NMOS is the sensitive diffusion, the node

voltage needs to be pulled down much farther to flip the output of the inverter. Based on this information, it would make more sense that the PMOS transistors would have higher sensitivity. However, looking at the physical layout again revealed that the sensitive PMOS drain nodes in the mirrored column orientation were 2x farther apart than the sensitive NMOS drain nodes, showing that the characteristic that may have the greatest effect of the likelihood of MCU/MBUs is the distance of sensitive nodes in the layout.

## 2.7.2   Digital Logic

### 2.7.2.1   Static Noise Margin

Fig. 2.22 shows the measured alpha particle soft error rates for each static noise margin (SNM) test where $V_{DD}$= 0.35V-0.5V. Two cases are shown in the figure: the top plot shows when the 'static_in' input is set to a logic '0', and the bottom plot shows when it is set to a logic '1', thereby reversing the effective SNM of each test. The smaller range of supply voltages is shown because larger sample sizes were accumulated at the lower voltages, which helped to even out the data trends. At supply voltages between 0.8V to 1.0V, only the lowest SNM (20%) test experienced any errors. From 0.6V to 0.8V, the 20% and 30% SNM tests both had errors. Once $V_{DD}$ was scaled down to 0.5V, the 50% SNM test began seeing errors. No errors ever occurred at any supply voltages with the SNM tests whose switching thresholds were greater than 50%.

It was found that the added wire capacitance from the automated place and route had a much more significant effect on the SER than originally expected. To adjust for this, error rates are first captured and then normalized for capacitive load using layout-extracted node capacitance. This allowed for a strict assessment of the effect of only the circuits under test on error rate, rather than any extra inadvertent layout effects. It should also be noted that, since adjusting the SNM requires an adjustment of P/N ratios, any interpretation of the results does need to take into account the change in diffusion area and parasitic capacitances of the transistors on top of the change in inverter switching threshold.
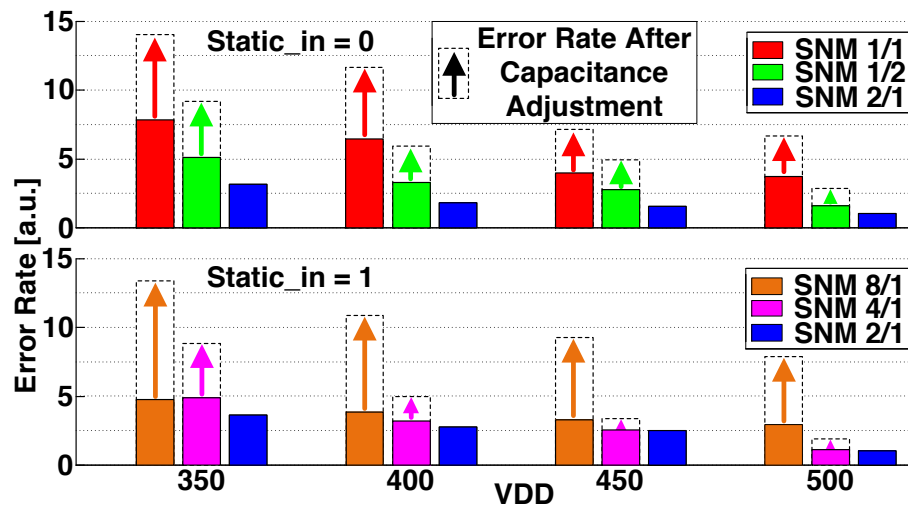


Figure 2.22: Static noise margin test error rates vs. $V_{DD}$

Fig. 2.22 shows that a decrease in SNM from 50% to 20% for the 'static_in'=0 case results in an increase in the unadjusted error rate by 3x at all supply voltages shown. Normalizing for the added wire capacitance (using the 50% SNM test as the baseline) increases the error rate difference. However, this likely leads to an

overestimation, as the parasitic capacitance of the 25% SNM test is larger (larger P/N ratios) and thus won't result in errors trends that match exactly with the difference in input threshold.

The 'static_in'=1 case showed slightly different error trends at all voltages, even resulting in the 35% SNM test having a higher unadjusted error rate than the 20% SNM test at $V_{DD}$=0.35V. This result was surprising, as the expectation is that the trends should mirror those of the 'static_in'=0 tests, since the tests themselves are mirrors of each other. A look back at the layout-extracted netlist shows that the wire capacitance difference due to automated layout leads to node capacitance values for the SNM 8/1 and SNM 4/1 tests that vary by as much as an order of magnitude from the SNM 1/1 and SNM 1/2 tests, respectively. In theory, these node capacitances should be the same, as the different tests consist of similarly sized inverters in reverse orientations (Fig. 2.9a). This observation is what prompted the calculation of the capacitance normalization in the first place, and once it was applied to these results, the data resembled trends that mirrored the 'static_in'=0 case, which was to be expected.

The most straightforward and significant conclusion to draw from the collected data is that a 50% SNM sees little-to-no errors until $V_{DD}$ scales down to the near-threshold region. The implications of this are significant, as combinational logic is typically designed with input thresholds near 50% to maintain equal propagation delay for all input transition possibilities. According to the results of these experiments, SETs are not a large issue for typical combinational circuits at nominal voltage. However, they become an issue that low-voltage designers have to seri-

ously address. If traditionally designed logic starts seeing a significant increase in SETs at lower-voltages, simply increasing the SNM of the same logic paths will prove detrimental to the already-decreased performance. Creative measures – such as altering the SNM of certain gates in some paths, utilizing only one transition of those gates to maintain reasonable timing constraints – will need to be taken. This, and other hardening measures, will come with overheads that low-voltage designers will need to consider.

## 2.7.2.2 NAND-based vs. NOR-based Logic

Fig. 2.23 shows the alpha particle experimental results of the NAND/NOR tests for $V_{DD}$=0.35V to 0.5V. The top plot compares the NAND test to the NOR-X test, while the bottom plot compares the NOR-X test to the NOR-Y test. Similar to the SNM plots, two data sets are shown for each test case, the strictly observed error rate and the data adjusted for layout-extracted node capacitance due to the larger-than-expected effect of wire capacitance.

All NAND/NOR chains saw 0 errors for supply voltages ranging from 1.0V down to 0.7V. At 0.7V, all NAND/NOR chains began to see errors. Interestingly enough, this was at a higher voltage than where the 50% SNM inverters began to see errors, implying that the NAND/NOR standard cells are more sensitive to SETs than the standard inverter. There are two reasons that can most likely explain this. First, the transistors in the NAND/NOR gates are larger, on average, than the transistor sizes in the SNM tests. And second, there are multiple transistors on

Figure 2.23: NAND and NOR test error rates vs. $V_{DD}$

the sensitive nodes. Both of these characteristics lead to increased collection area, increasing the upset probability at any supply voltage.

A comparison between the NAND and the NOR-X tests focus on the difference in collection area and output node capacitance. The NAND test shows 1.5x higher error rates at all voltages before any capacitance normalization, showing that the increased collection area from the larger devices on the sensitive node in the NAND may have a larger effect than the added parasitic capacitance. In fact, the data after capacitance adjustment shows a larger error rate increase for the NOR-X test than the NAND test, which means that the added wire capacitance on the NOR-X test was large enough to affect the initial measurement results.

A comparison between the NOR-X and NOR-Y tests shows that the soft error rate can increase based upon which input is set low, due to multiple sensitive drain nodes within a single logic gate, as described in Fig. 2.10. The non-normalized

results show that the NOR-Y test has a 1.25x higher error rate than the NOR-X test. This implies that having multiple sensitive nodes in some gates throughout the chain can increase the susceptibility to soft errors.

For low-voltage designs, the measurement results show that it is beneficial to implement NOR-based logic over NAND-based logic, as the smaller collection area of the transistors in the NOR-gates help to decrease SET probability. This also has the added benefit of lower area per cell. The NOR-X vs. NOR-Y test also shows that what state the static logic gate is held at can have a small effect on the error rate. Careful implementation can be utilized so that fewer nodes are exposed in static logic, leading to more reliable combinational logic designs.

### 2.7.2.3   Inverter Size/Chain Length

Error rate vs. inverter chain length under alpha particle radiation is shown in Fig. 2.24. Three separate lines are plotted for the 1X, 4X, and 8X sized inverter chains at $V_{DD}$=0.35V. The 0.35V measurements were chosen as the supply voltage to plot because the 8X inverter chain did not have any errors at supply voltages higher than 0.4V, with a sufficient sample size only at 0.35V. For chains consisting of inverters sized at 1X and 4X, errors continued to increase steadily as the chain length increased all of the way up to the 32-length chain, the largest included on the test chip. As previously described, the assumption is that as long as the error rate increases with increasing chain length that transient errors could continue to propagate through a chain length as long 32 gates. Ideally, the plot would have

Figure 2.24: Chain length vs. error rate for 3 different inverter sizes

a linear increase. However, the results are a bit uneven, as the increasing slope never really remains constant. This is probably due to both too few data points across chain lengths as well as small sample size noise. The 8X chain shows a different trend than the 1X and 4X curves, however, as the error rate plateaus after a chain length of 8 – even decreasing a bit at data points for the longer chain lengths (though they remain within each others error bars) – leading to the conclusion that transient errors could not propagate through a length of more than 8 inverters.

The results from this experiment are what would be expected, as large, slow devices should limit the length of error propagation. What's more interesting is

that at the shorter chain lengths (2, 4, and 8), the INV8X chain error rate was roughly similar to the INV1X and INV4X chains. This shows that the larger sized inverters, and the larger capacitance that comes with them, are not much more resistant to the initial existence of SETs, most likely because of their increased diffusion area. Still, because low-voltage designs are typically not designed for the fastest speeds, using larger devices on non-critical computational paths as well as on clock trees to prevent glitches is suggested. Additionally, the use of larger sized devices will also have the added benefit of process variation mitigation.

Fig. 2.25 shows a scatter plot of all logic test errors rates vs. $V_{DD}$ under alpha irradiation. Only the test with the smallest static noise margin exhibited any errors at 1.0V, and at that voltage, the error rate was very small (1 error every 2.5 hours of irradiation, resulting in an inverter cross-section 1/10th of the 6T SRAM). The error rate for this test increases by 6.45x from 1.0V-0.35V. A drop is seen in the error rate for this test when decreasing the supply voltage from 1.0V to 0.9V, though this is most likely due to small error sample size.

This figure provides useful information in that it shows that standard combinational logic sensitivity significantly increases once $V_{DD}$ drops below 0.6V. As previously mentioned, the only tests where errors occurred above $V_{DD}$=0.7V were the SNM tests with the 2 smallest switching thresholds (20% and 30%). Logic of this type is not traditionally used in current designs, so the trends from the other tests (which use standard cells, and typical place and route techniques) are more indicative of what to expect in actual designs. The NAND/NOR tests, which consist of standard cells, do not have any errors until 0.7V and see error rate in-

Figure 2.25: All logic test error rates vs. $V_{DD}$

creases ranging from 3x to 5x as $V_{DD}$ is lowered from 0.7V to 0.35V. The inverter size/chain length tests, which also use standard cells, do not experience any errors until $V_{DD}$=0.5V, and see error rate increases ranging from 2x to 6.5x as $V_{DD}$ is scaled from 0.5V to 0.35V.

The overall voltage scaling trends indicate that combinational logic SER will be an increasing issue that needs to be addressed in sub/near-threshold designs, especially for reliability-critical applications. What's interesting to note is that for the tests utilizing standard gates and layout techniques, no errors occurred at nominal voltage over 1 week of accelerated alpha experiments. While this does not necessarily mean that radiation-induced combinational logic errors do not oc-

cur in standard 65nm designs, the significant increase in measured error rate at $V_{DD}$=0.5V and below is a strong indicator that logic SER may end up dominating chip-level SER in low-voltage designs.

## 2.8   Design analysis

This chapter detailed the design and implementation of the first version of a soft error test reference platform. This initial implementation was done in a TSMC 65nm CMOS process, with plans to port the design to more advanced processes in the future. A single test-chip working within the entire framework of the design was verified using an Am-241 alpha point source provided by Oregon State Universitys Radiation Center. The entire setup was implemented using the neutron beam in ICE House I at LANSCE.

Many key conclusions were taken from the data resulting from experiments run on this first version. Neutron-induced 6T memory SER showed a 6.45x increase when supply voltage scaled from 1.0V down to 0.33V, with a distinct increase in slope once the threshold voltage ( 0.35V) was approached. 8T data trends were similar with just slightly lower SER in most cases, which is to be expected. MBUs were also shown to be largely affected by layout techniques, and with increasing MBU rates at low supply voltages, (including more common occurrences of TBUs) more SEU aware memory layout techniques may prove to be highly beneficial.

Data from alpha particle combinational logic experiments was also collected. The results showed that, overall, logic SER will become more of a problem at

low supply voltages. Specifically, typical inverter static noise margins optimized for best performance were sufficient to prevent SETs at nominal supply voltage, according to alpha radiation test results. However this 50% SNM begins to experience more errors at lower voltages, prompting the need for alternative methods for SET mitigation, as simply adjusting the SNM will hurt the circuit functionality and performance. NAND/NOR tests also exposed information on how individual circuit differences and tendencies can affect SER. And, lastly, inverter size/chain length tests helped to discover how SET propagation can be limited at low supply voltages.

Based on the data that was collected for both the SRAM and logic test structures, preliminary guidelines can be made for effective circuit sizing for future process node low-voltage SER mitigation. Increasing the SNM of the SRAM cells to match the 50% threshold of standard logic will improve the SER of SRAM with low-voltage retention states, with the drawback of decreasing the SRAM density and performance. Additionally, for digital logic, larger cells for decreased SET propagation combined with increased logic depth will decrease the probability of an SET being clocked into a sequential cell and manifesting itself as an error. These are reasonable changes for combinational logic in low-voltage designs, as target applications for near-threshold operation will have lower performance requirements.

Being that the main goal of this reference design is to develop a platform that is easily scalable across both voltage and process, it is important to treat each implementation as an opportunity to improve the platform for future generations. While this version of the test setup was mostly successful, there are some aspects

of it that will need to be improved, starting with the chip-level design. For the memory tests, the first area for improvement is the use of realistically sized high density 6T and 8T bitcells. The benefits of this are twofold, as not only will 2x more cells be able to fit on the same die area, but more applicable and realistic MBU measurements can be taken. This change can be made quite easily, as the cells can be custom laid-out using SRAM design rules. The key will be getting clearance from the foundry a sufficient amount of time before the tapeout deadline. Including the interruptible 10T bitcell on future test chips will also require more careful design, as well as the use of more dense SRAM layout rules.

Many potential improvements for the logic test structure also appeared during testing. Throughout both the alpha and neutron experiments, the error counts never came near the maximum storage value of 32 of the 5-bit counter. For future versions, the counter size could be reduced to as small as 3-bits. The TMR on the logic scan chain also could be removed, as no errors were found during the duration of any experiments. Both of these changes will free up a significant amount of area and allow us to increase the cross-section per test. This is very important, as the data collection efficiency needs to increase a large amount to acquire a sufficient amount of data during limited neutron beam test time. Finally, our efforts to assess the effect of different circuit characteristics on logic SER were complicated by the added wire capacitance from standard cell place and route layout effects. The added capacitance from these wires had a unexpectedly large effect on total node capacitance, and therefore affected the final error counts and skewed the data to the point where it resembled unexpected trends. This issue can be fixed by using

more careful custom layout of the digital logic tests, to limit the wire capacitance overhead.

The final presented logic error rates were only for the asynchronous logic tests, as the off-chip PLL on the master board was not debugged successfully before tests needed to be run. While the error rates for these asynchronous tests are useful, the chances that many of the observed SETs actually get clocked into a sequential element (thus establishing themselves as errors) is highly dependent on the clock frequency (timing masking). This is a key factor into the full understanding of SETs in combinational logic and will need to be added for future versions of the test platform. The results from the first version, as they stand, still provide valuable information as they give important details about the sensitivities of the circuit structures themselves.

Beyond the chip-level portion of the test platform many issues were found during board design that should be improved upon, mostly with the purpose of simplifying the assembly process and test setup. While fitting 8 chips (each with 72 pads) into a circle with a 2-inch diameter to accommodate the beam path was initially successful, the complexity of the board design and having to wire bond the chips directly on board significantly lowered the yield of the chips, as only 24 out of the 80 chips were test-able. Additionally, using chip-on-board wire bonding was very expensive. Redesigning the sub-boards and allowing chips to be packaged will cut down on cost and complexity. This can be approached either by placing fewer chips on a board and adding more sub-boards to the designs, or by placing chips on both sides of each board. The second method is preferred for lower cost,

and the ability use a similar master board design.

For the master board design, power distribution became an issue as more test chips operated off of the same four shared supplies. During neutron beam testing, it was necessary to use the sense inputs on the power supplies to stabilize the output voltage levels. For future versions, it will be necessary to partition the shared test chip supply domains into smaller subsets while increasing the total number of power supplies used.

Overall, the first version of the reference design was very successful, both for data collection and for learning how to construct proper tests for ideal SER characterization. With continued improvements in future test chip versions allowing for detailed and specific data collection, the development of simulation models should be achievable, providing future low-voltage circuit engineers with a valuable tool to make their designs more efficient and reliable.

# Chapter 3: Synctium-I: A 10-Lane, Near-Threshold SIMD Processor Incorporating Timing Variation Resiliency Techniques

---

## 3.1 Introduction

An increasing number of computing applications are requiring higher performance per watt. For many applications, data-parallel processing architectures can be exploited to deliver both high performance and low power consumption. For example, in the domain of biomedical sensors, EEG artifact separation requires a throughput of 0.016 - 1.9 GOPS [45]. As biomedical sensors are increasingly adopted in medical implants and future body-area networks that are energy constrained (such as battery-powered or energy-harvested environments), the lowest power consumed while satisfying a real-time signal-processing throughput will be necessary.

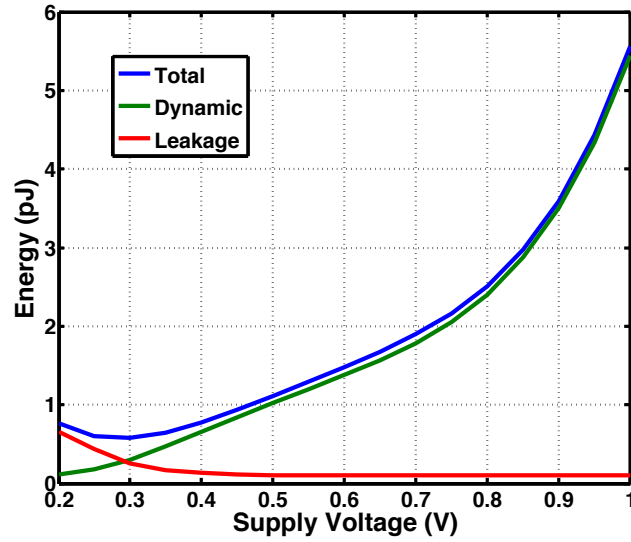Other data parallel applications, such as video processing within a mobile cellphone or a wireless camera will also require high computational throughput. For example H.264 video encoding and decoding with 30fps @ 1Mbps, require approximately 1.0-2.5 GOPS and 0.2-1.2 GOPS, respectively [46]. Overall, because battery lifetime and charging overhead are important design considerations, the goal is to

design a programmable processor that can simultaneously achieve massive parallel processing throughput at the lowest possible energy/computation.

### 3.1.1  Near-Threshold Operation

A popular method for lowering processor power is to operate the supply in the sub/near-threshold voltage (NTV) region [1], [47]. Because of the square-law dependence on supply voltage, operation at near-threshold can result in an energy/computation decrease of approximately 5-10X versus conventional super-threshold operation. Operation at sub-threshold (where the minimum energy point exists) can yield an energy/computation decrease of 20X [1], at the cost of significantly degraded delay. For example, a sub-threshold, 16-bit 1024-point FFT test-chip operating at $V_{DD}$=0.35V (with $V_{TH} = 0.45$V) achieves 155nJ per FFT at 10KHz [48], while an 8-bit processor achieves 540 fJ / operation for $V_{DD}$=0.3V operating at 160KHz [49].

While operation at sub/near-threshold voltages can yield a significant energy savings, it also introduces two problems: severe performance degradation, and an increase in timing variation due to heightened sensitivity to process variations. Timing delay spread (the difference in minimum to maximum delay) elevates due to both static (random dopant fluctuations) and dynamic (temperature, supply voltage droop, input vector dependent, and radiation-induced soft errors) variations at low supply voltages, and will accelerate with continued decrease in minimum feature size [50].

(a)



(b)

Figure 3.1: Simulations of a 16-bit fixed point multiplier in 45nm SOI: (a) static and dynamic energy vs. $V_{DD}$; (b) delay vs. $V_{DD}$

For example, [1] reports a 10X performance degradation when the supply voltage is lowered into NTV, and as much as a 200X performance decrease when

operating in sub-threshold. Fig. 3.1 shows the results of our own simulations of a 16-bit fixed-point multiplier in 45nm SOI across scaled supply voltage. The energy relationship is shown in Fig. 3.1a, where the dynamic energy dominates the total energy consumption until just before the minimum energy point is reached. A 5X energy decrease occurs as the supply voltage is lowered from 1.0V to 0.5V ($V_{TH} = 0.37V$). At the minimum energy point ($V_{DD} = 0.3V$), the total energy consumption decreases by 10X. The plot of delay vs. $V_{DD}$ (Fig. 3.1b) shows a 4.5X performance decrease when the supply voltage is lowered from 1.0V to 0.5V, and a 30X degradation when $V_{DD} = 0.3V$. This severe performance decrease suggests that, despite the slightly improved energy-efficiency of sub-threshold operation, near-threshold (NTV) operation is a more reasonable option for throughput-constrained low-power applications.

Fig. 3.2a shows the effect of supply voltage scaling on the timing delay of a 16-bit fixed-point multiplier in 45nm-SOI. Each box-plot represents an identical set of 50 random input vectors simulated at supply voltages ranging from $V_{DD}$ = 1.0V to $V_{DD}$ = 0.5V at 100mV increments. Lowering $V_{DD}$ from 1.0V down to 0.5V results in a min-max delay variation increase of 4.5X. Fig. 3.2b shows the effect of spatial variation on the clock frequency of multiple identical units placed on the same die. Randomly seeded Monte Carlo simulations were run for a set of 50 random input vectors for 10 identical 16-bit multipliers at $V_{DD}$ = 0.5V. As observed, the maximum delay of multiplier-4 is 25% slower than the maximum delay of multiplier-8. It can also be noted that the distribution of delays within the same multiplier varies significantly for each case. Hence, while near-

Figure 3.2: Simulated delay variation for a 16-bit fixed point multiplier for: (a) Scaled supply voltage; (b) Spatial variation of 10 multipliers on the same die

threshold operation shows a large potential for energy-improvement, the inability to predict delay and achieve timing closure has prevented its widespread adoption in commercial products.

## 3.1.2    Error Detection and Correction

The most basic approach to dealing with variation at nominal supply voltage is to operate with timing guard bands, and set the clock frequency safely below the worst case path to where any possible timing variations do not cause an error. At lower supply voltages, with the increase in delay min-max variation, the timing guard bands need to be increased. With the clock frequency already significantly

decreasing at lower supply voltages, operating with larger guard bands becomes impractical, and an alternate solution becomes necessary. This subsection explores different possible solutions that have been developed. The feasibility of each idea will be discussed, with a focus on any potential issues that could arise while implementing the idea at low-voltages.

Prior work [48], [49] has explored a number of more straightforward solutions to combat the effects of variation. Using non-minim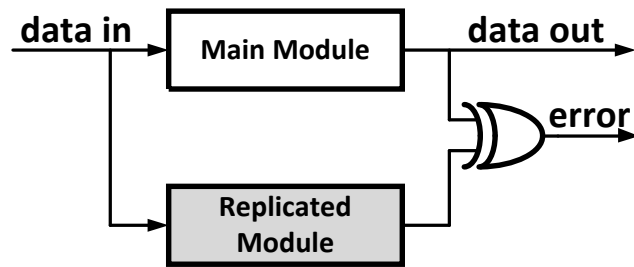al device sizes or biasing the body connection of the transistors (using the body effect to adjust the transistor threshold voltage) can protect against an increased number of radiation-induced soft error current pulses, as well as supply voltage droop. However, employing larger transistor sizes diminishes the benefits of technology scaling, and body biasing requires increased testing time, calibration, and requires the use of additional voltage domains, which in turn incurs larger overheads for additional voltage regulators. Longer logic chains are also presented as a possibility, as the increased logic depth averages out the effect of variation on the logic chain. Increased logic depth will also increase area and power overhead, as well as potentially increasing the worst-case delay, resulting in lower clock frequency. Beyond the issues presented, these solutions may be unable to completely address the increased timing uncertainties expected with continued deep sub-micron processes.

Other approaches have explored various types of circuit and architectural-level implementation to perform error detection. At the architectural level, Dual Modular Redundancy (DMR) serves as the simplest form of error-protection (Fig. 3.3a). In this method, the protected module is replicated such that both the original

module and its copy share their input signals. Outputs are compared and a mismatch indicates an error. In the event of an error, the system has to restore its last verified stat and re-execute from that point. While having the advantage of design simplicity in using exact replication, this approach's disadvantages include high area and power/energy overheads, as the design sized is more than doubled. In addition it requires an alternate error-correction mechanism that can restore the last verified state and re-execute, which will create even more area/power overheads, as well as throughput loss with the extra clock cycles of recovery. Triple



(a)



(b)

Figure 3.3: Architecture level error-protection techniques (a) Dual Modular Redundancy; (b) Triple Modular Redundancy

Modular Redundancy (TMR) is similar to DMR, but this approach utilizes two

replicated instances in addition to the original module (Fig. 3.3b). The outputs of the three instances are compared and a majority voter mechanism is used to select the outcome signals. Based on the assumption that only one of the instances (at most) will fail at any given point, this approach eliminates the need for the roll-back and replay mechanism that was required for DMR, allowing for instant recovery. However, adding a third redundant instance of the module further increases the area and power overhead. Also, the critical path is increased, as the computation and voting mechanism must occur within the same clock cycle.

Error detection approaches at the circuit level require more meticulous design – focusing on the computational paths that are most likely to produce an error – while providing lower area/power overheads when compared with modular redundancy. One circuit-level approach is the use of Tunable Replica Circuits [51]. These circuits are composed of a number of digital cells, such as inverters, NAND, NOR, adders, and metal wires that are tunable to a given delay time (Fig. 3.4). The replicas are affected by process variations and aging in a similar way to the critical path. Once the replica is tuned to the critical path, it will replicate the path delay as it changes due to these variations. The TRCs can be used to report the critical path delay using a thermometer code or perform dynamic error detection. TRCs are able to detect errors without introducing additional components or time delays to the data path. After being tuned once, the TRCs will mirror slightly worse than the critical delay path to ensure that any timing violation will be observed. However, TRCs do suffer from a few drawbacks. Because the circuit is only a worst case replica, it is possible that they will trigger error responses when
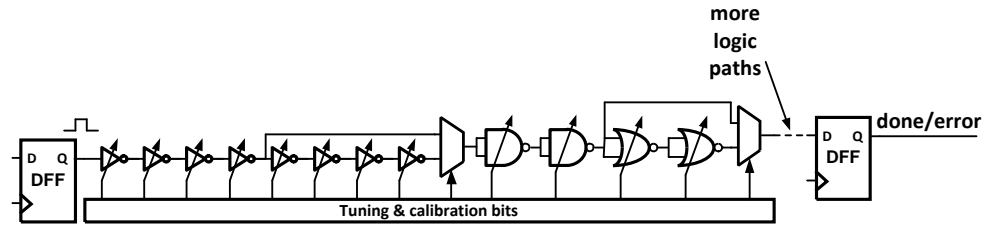
Figure 3.4: A typical Tunable Replica Circuit (TRC)

there was no timing violation in the actual data path. The circuits themselves also take up area and power. Finally, TRCs cannot adapt to unique delay paths, only a simple worst case. This can make them hard to calibrate correctly under extreme variations as the TRC's timing margin needs to be large enough to guarantee all errors will be caught correctly.

Other circuit level timing error detection methods to combat dynamic variations take the approach of placing the error detection circuitry on the worst case computational paths [52]. One such approach is Razor-I [53]. This circuit was introduced to detect errors by double sampling the output data from the unit under test (Fig. 3.5a). This is done by implementing a positive level-sensitive error-path latch that is clocked on a delayed clock. Fig. 3.5b shows the timing diagram of the Razor-I error detection circuit. The delayed clock is set in phase with the data-path clock, but at half of the duty cycle. As long as the data arrives at the Razor-I flip-flop before the rising clock edge (without violating setup-time constraints), the outputs of the the data-path flip-flop and error-path latch will be the same, and no error will be flagged. If the data arrives after the clock edge but within the positive phase of the delayed clock, the outputs of the error-path

latch and data-path flip-flop may be different. In this case, an error signal is set and a recovery method is activated to fix the incorrect data resulting from the timing error. While Razor-I is an effective method of error detection, there are



(a)



(b)

Figure 3.5: Razor error detection sequential circuits (a) Circuit schematic; (b) Timing diagram

a few requirements to ensure proper functionality. If the arriving data violates setup-time constraints, the data-path flip-flop can become metastable. Razor-I uses extra circuitry to determine if the flop-flop could become metastable. If so, it is treated as an error and appropriately corrected. There are also limitations on the size of the timing window for error detection, as it must be smaller than the

fastest path through the logic. If it violates this constraint, a race condition from the next computation will introduce a false error. To allow for a larger detection window, buffers are inserted into the fastest paths, resulting in increased area and power overheads beyond just the Razor circuits themselves. These overheads are minimized during super-threshold operation by placing Razor circuits on only the few slowest paths and buffering only the few fastest paths to maintain a beneficial detection window. Unfortunately, operating Razor in near-threshold can incur even larger overheads as the increased timing uncertainty requires Razor flip-flops on more output bits, and larger min-max delay spread (as shown in Fig. 3.2a) requires buffer insertion on additional fast paths of the logic.

Once an error is detected through the use of circuit or architectural-level error detection techniques, a recovery procedure is enabled to handle the incorrect data that results from the timing error and return the pipeline to its correct state. Clock gating [54] is the most straight forward recovery method (Fig. 3.6), as it simply stalls the entire pipeline in the event of an error, either allowing the computation an extra clock cycle to fully complete, or for the data to correct itself using the value stored in the shadow latch of the Razor flip-flop. While it requires few architectural changes and little logic overhead, this method can put a serious limitation on the clock frequency because it requires the stall signal to be propagated to all of the pipeline stages in a single cycle. This can present a significant problem in larger and more complex designs, where it may take several clock cycles just to propagate the clock signal through a clock distribution network, which cannot be halted in only one cycle.

Figure 3.6: (a) Pipeline modification for Clock Gating error recovery method; (b) Clock Gating pipeline data path with errors.

Other methods, such as Micro-Rollback [55] remove the need for a global stall signal. Micro-Rollback is a technique that saves a queue of previous instructions and operands at each pipeline stage. After a successful operation, each stage saves the results as they are passed to the next stage. If an error is detected, instructions can be restarted at their last known correct state in the pipeline. Once an error is detected, a signal is propagated with the instruction as is continues through the pipeline. Once the instruction reaches the write-back stage, rollback logic is activated to stop the instruction from being written and sends a signal to all of the pipeline stages to roll back and retry the instruction. The Micro-Rollback architecture is illustrated in (Fig. 3.7a). In its original design, instructions were simply replayed with the hope of the error being resolved at a later time. This

is not very robust but with small architectural modifications this method can be redesigned to replay instructions twice or three times to ensure no error is produced a second time.



(a)



(b)

Figure 3.7: (a) Pipeline modification for Micro-Rollback error recovery method; (b) Micro-Rollback pipeline data path with errors.

The elimination of a global stall signal makes techniques such as Micro-Rollback useful for designs where traditional clock gating is not practical. However, it does introduce significantly more logic overhead. Additionally, the recovery penalty

for this method, depending on pipeline depth and where in the pipeline the error occurred, can be several cycles long – much larger than the one-cycle penalty for clock gating.

### 3.1.3 Error Resiliency on a SIMD Pipeline

Utilizing highly parallel architectures with lower clock frequencies under near-threshold operation can help achieve the required performance for many low-power applications. In order to maximize the benefit of energy efficiency attained with near-threshold operation, programmable architectures must exhibit efficient control, such as with wide SIMD/vector execution [56], [57]. Unfortunately, wide SIMD architectures with a large number of ALUs exacerbates the worst-case timing closure problem, which is further degraded by operation in NTV. As shown in Fig. 3.2 and in [50], operation in near-threshold can result in a significant degradation in logic delay across Monte Carlo variation. Even worse, because transistor variations such as random-dopant fluctuation are not related to any spatial distances, parallel functional units will not exhibit spatial correlations that would enable similar timing delays between nearby processing cores [58]. Operating these wide SIMD architectures at near-threshold with sufficient timing guard bands can potentially result in performance losses so large that it negates the increased throughput of the parallel architecture compared with a single lane, rendering this design choice impractical.

Implementing Razor on a SIMD pipeline will enable timing speculation and er-

ror detection on many parallel cores, helping to recover much of the lost throughput of near-threshold operation. Unfortunately, because the inherent nature of SIMD is for all of the lanes to operate synchronously in lockstep on a global clock, a couple of problems can arise during the error recovery process. The first is that the use of a wide-SIMD architecture eliminates the possibility of using traditional clock gating, as the propagation of a global stall signal within one clock cycle will degrade the performance to an unacceptable level. The second problem is that any error that occurs in a single stage of any core will cause a global stall of all the parallel cores. This will have a significant effect on throughput, especially for larger designs, as the effect of individual errors on performance is multiplied linearly with the number of SIMD lanes. This overhead is compounded by the fact that the probability of an error occurring in any one stage of any lane is effectively larger with a parallel pipeline because the chance of a timing violation in one functional unit is independent from the other functional units that are processing different inputs (Fig. 3.2b). Addressing the first problem by replacing clock gating with an alternative recovery technique (such as Micro-Rollback) will essentially eliminate the need for a global stall signal. However, there will be a very large total area overhead, as extra logic will be added to all lanes, and the throughput loss will be large, as the number of cycles delay for each error will be multiplied by the total number of lanes in the design.

Fig. 3.8 shows the fraction of peak throughput versus the presence of timing errors for different architectures. Because all of the SIMD lanes will need to stall for any error within a single lane, a wide parallel architecture will lose significant

Figure 3.8: Throughput degradation of SIMD architectures of different widths in the presence of timing errors, versus a scalar pipeline.

performance even at relatively low error rates. For example, at a 1% error rate, a 32-wide SIMD architecture will see a throughput decrease greater than 43% when compared with its peak throughput, versus a scalar pipeline. It should also be noted that while not directly shown in the figure, if clock gating is the recovery method of choice the peak throughput will not increase linearly as the number of SIMD lanes increases, due to the limitation of the global stall signal propagation over an increasing distance on clock frequency. As a result of this significant throughput degradation, the effectiveness of timing speculation for SIMD architectures is essentially negligible, as the throughput loss is comparable to conventional operation of the processor at a frequency significantly below the guard-bands.

This chapter will detail an implementation of variation resiliency techniques

intended to further enable near-threshold operation of a SIMD parallel architecture, as first described in [59], and then implemented on-chip and tested in [60]. The rest of this chapter is organized as follows. The two proposed methods for variation resiliency: the Decoupled Parallel SIMD Pipeline (DPSP) and Pipeline Weaving, are described in Section 2.2. Section 2.3 details the micro-architecture of the test-chip, specifically highlighting the physical implementation of our proposed methods. The final chip design, system design, and test methodology will be described in section 2.4. Silicon measurement results for the fabricated chip will be given in Section 2.5. Finally, Section 2.6 concludes the chapter with an overview of the project as well as identification of areas of improvement in the design.

## 3.2    Variation Resiliency Methods

### 3.2.1    Decoupled Parallel SIMD Pipeline

The Decoupled Parallel SIMD Pipeline (DPSP) [59] can be implemented to combat the effects of dynamic variations on a SIMD architecture. With DPSP, all functional units in the SIMD organization still execute the same instructions in the same order, but the lockstep operation of the parallel pipelines is no longer required, thereby allowing each pipeline to deal with errors independently. Because a global stall signal is no longer required, clock gating because a possibility, as stall signals kept local do not put a limitation on the clock frequency. The DPSP utilizes Decoupling Queues (DQs) at the beginning of each SIMD lane, which con-

tinuously pass instructions sent from the sequencer to the individual lanes. If an error is detected within a single lane, only that lane is stalled and proceeds with error recovery. All other lanes proceed normally, as if no error had occurred. While that individual lane stalls, the next upcoming instructions will fill into that lanes DQ. In the event that a lanes DQ is full, the entire pipeline will have to stall, in order to free up the full DQ with a vacancy. In this case, a full signal is sent from the full DQ to the sequencer, still avoiding the propagation of a signal over long distances.

Due to this parallel lane decoupling, proper synchronization across all lanes will require the introduction of micro-barriers. Instructions that enable the movement of data between lanes (shuffles) or require memory accesses (loads and stores) cannot be decoupled, as there are dependencies between lanes. These types of operations must be executed in the proper order – each lane in a shuffle must wait for its producer lane to supply the correct value, and with a memory load/store the data must be read/written from memory as the program dictates. To address these issues, micro-barriers are utilized to synchronize the lanes between the decoupled parallel SIMD pipelines. Once the micro-barrier instruction arrives at the front of the DQ, it prevents the DQ from proceeding on to the next instruction, until all of the other DQs are aligned in time at that same micro-barrier instruction. At this point, all of the lanes can execute and again operate in lock step synchronicity.

Fig. 3.9 shows the potential throughput improvement that results from the proposed DPSP design. Comparing the throughput of the proposed DPSP scheme with a conventional SIMD scheme, the DPSP shows a significant improvement over

Figure 3.9: Throughput improvement of DPSP for 16 and 32-wide SIMD architectures in the presence of timing errors.

SIMD for the same error rate. For example, it is observed that for an error rate of 1%, the fraction of peak improvement for a 32-wide DPSP is 39% more than a conventional SIMD without decoupling, with both architectures utilizing a Razor-like error detection method. It should be noted once again that, assuming that clock gating is the error recovery method of choice, the peak throughput for SIMD versus DPSP will be much different. Because the stall signal is kept local for DPSP, the clock frequency will be much higher for DPSP compared to regular SIMD for any number of lanes. The exact relationship is dependent on the physical layout of each of the designs. However, it can be assumed that for larger numbers of lanes, the max clock frequency difference between DPSP and SIMD will increase.

### 3.2.2 Pipeline Weaving

While DPSP can combat the effects of dynamic timing variations (i.e. temperature, supply droop, soft errors), other methods are necessary to combat any static timing variation. Process variations can create large delay variations between the various stages of a SIMD pipeline, degrading both the energy-efficiency and performance of the processor. For example, if the timing variation between the different ALUs is large, both performance and energy-efficiency will be impacted, as the system clock will be set to the worst-case slowest paths, and the faster blocks with high leakage will be clocked with longer cycle times, increasing the integrated leakage energy.

Pipeline weaving [59] is proposed to allow fine-grained sparing of multiple blocks within the parallel pipeline. To implement pipeline weaving, a small number of redundant components are added to the SIMD pipeline. The optimal number of redundant components (2, for this implementation) was determined by the delay spread observed through Monte Carlo timing simulations. Selecting the proper number is important, as there is a large area tradeoff for adding extra blocks (no power overhead, as blocks can be power gated). Adding too many redundant components will result in too large of an area overhead with minimal frequency increase per block added. Too few extra components would not result in enough coverage to deal with the static variation existing in all of the functional lanes, and would not sufficiently increase the max clock frequency.

The basic conceptual implementation of pipeline weaving for error resiliency is

Figure 3.10: Lane weaving connections for a 10-lane (8+2) SIMD architecture

shown in Fig. 3.10. Within each lane, each pipeline component is connected to the two neighboring lane's components in the subsequent pipeline stage [61], [62]. Blocks that are initially characterized to be the slowest functioning within its pipeline stage (Marked as a red X in Fig. 3.10), are bypassed from the active paths. While weaving is possible with only two connections (top and bottom adjacent lanes), a third connection is made directly to the next pipeline stage in the same lane. This 3-port weaving provides greater flexibility for the weave orientation, cutting down on the average wire length connectivity, as connections

within the same lane are shorter than those to adjacent lanes. The observation can be made that the top-most and bottom-most lanes have only two possible connections from each stage to the next. It is possible that if, for example, the bottom two RF blocks in Fig. 3.10 are the slowest functioning, the bottom DQ will have no available connections. In this case, the bottom DQ will need to be deactivated, and a potentially slower functioning DQ will take its place.

## 3.3  Chip Implementation

This section discusses the chip implementation of the DPSP and lane weaving on the parallel SIMD pipeline. It will first discuss the details of the decoupling queues, the queue depth chosen and what signals are necessary for proper operation. Then the entire pipeline will be shown, with the error detection and recovery logic highlighted. The full processor pipeline will then be discussed, with the lane weaving implementation explained, and the most complex component, the ALU, will be discussed in detail.

Fig. 3.11 shows the implementation of the decoupling queues within our test chip. During each clock cycle, the instruction queue (which can be pre-programmed with up to 128 instructions) broadcasts a single instruction to all lanes. Each lane incorporates an eight-deep decoupling queue that is controlled by an isolated stall signal from its lane only. The depth of the queue was chosen based on our expectation of error rate resulting from monte carlo simulations across supply voltage. This decision is important, as the tradeoff is increased area vs. throughput

given higher error rates, i.e. a larger queue will allow for more errors to occur before needing to stall all of the lanes, but at low error rates will only incur area overheads, and may not come close to being full.



Figure 3.11: Test-chip implementation of the sequencer and decoupling queues.

In the event of a stall the queue will cease to output new instructions, while still accepting instructions from the sequencer and placing them in the vacant slots. Once any individual queue is full, it activates a control signal that is sent to the IQ control block, stalling the sequencer from providing any new instructions until a vacant slot opens in the full DQ. If a micro-barrier instruction is executed, each decoupling queue sends a Bar signal to the IQ control block when the micro-barrier instruction reaches the front of the queue. In response, the IQ control sends a Barstall signal back to the DQ to ensure that it waits until all of the queues have

synchronized. Once all of the DQs have realigned the micro-barrier instruction to the front of the DQ, the Barstall signal is set low, and normal queue operation restarts with all of the lanes operating in lockstep.

An example case of the different states that each of the different DQ's can be in is shown in Fig. 3.11. Lane 0 and lane 1 have both experienced 5 and 7 errors in their pipelines, respectively. If one more error was to occur in lane 1, that DQ would send out the full signal to the IQ control to stall the IQ. Lane 9 has seen comparatively fewer errors than either lanes 0 or 1, however lane 9 has a micro-barrier instruction (in this case, a load) at the front of its DQ. Because of this, lane 9 will stall for at least 5 clock cycles (assuming that the lane 1 DQ has the largest number of instructions in it queue) to allow itself to sync with all of the other lanes.

The detailed implementation of error detection and recovery is shown in Fig. 3.12 . Razor-I circuits were used at the output of four different pipeline stages: the decoupling queues, register files, and two cycles within the ALU. Because the target of the pipeline was for near-threshold operation, Razor circuits were placed on every output bit of all four pipeline stages, not just the worst-case paths. Minimum path buffer insertion was performed during synthesis, taking extra precautions to ensure every path was balanced to allow for the maximum timing detection window under the increase variation resulting from low-voltage operation. This incurred significant area and power penalties, but ensured computational accuracy, which is especially important given the increased timing variability expected with near-threshold operation. Due to the use of decoupling queues, stall signals were kept

local within each active lane. This, combined with the relative simplicity of our pipeline and the lower clock speeds at near-threshold, allowed us to perform error recovery using clock gating with only a one-cycle penalty.
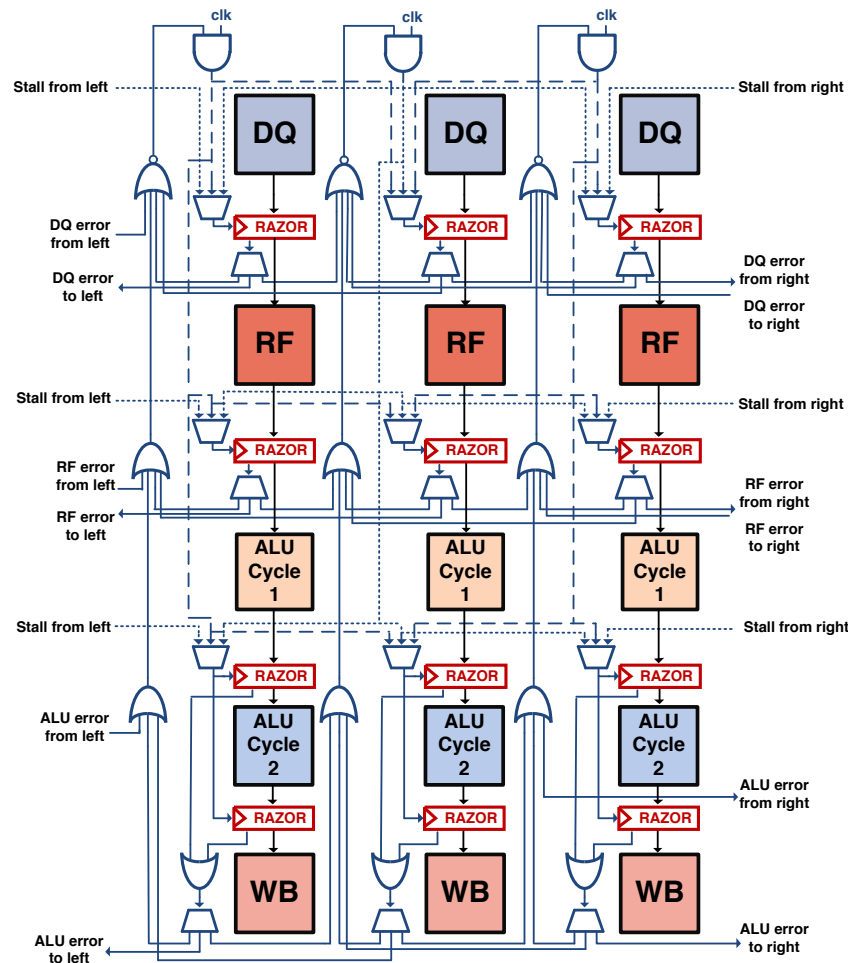


Figure 3.12: Implementation of error detection/recovery on the pipeline.

Lane weaving was implemented using 3:1 multiplexors at the beginning of three different pipeline stages: the register files, the ALU, and the write-back stage. In

order to identify the most effective weave orientation, delay characterization of all the units in the design was first performed. This characterization was started by using the STA worst-case input vectors that were found from timing analysis. Next, a number of input vectors that stressed computational paths close to the worst-case path were tested, in addition to a set of randomly generated inputs. For each set of inputs, the values were scanned into the beginning of the pipeline stage, clocked once, and then scanned out from the end of the pipeline stage. The clock frequency was increased until an incorrect value was scanned out. Once the two slowest-functioning blocks in each pipeline stage are identified, control bits are scanned into a separate weave control block to configure the desired weave orientation. While this calibration sequence can be time consuming, all of the blocks can be calibrated concurrently, cutting back on initialization time.

One issue that arrises with incorporating lane weaving in conjunction with error detection/correction is that active lanes may not simply follow straight through the pipeline, rather. Because of this it is necessary to ensure that the stall signals remain within each 'active' lane, no matter what physical path it takes through the pipeline. To accomplish this, extra circuitry was added to the stall signal paths. The select input for the multiplexors controlling the path of the stall signals is the same as the weave control input which is scanned in at startup, this ensures that the stall signal will follow the path of the data through each active lane. Based on simulation, the additional delay from the multiplexors added to the stall signal paths were found not to have a detrimental effect on the max clock frequency.

Fig. 3.13 shows lanes 0, 1, and 9 of the Synctium architecture. Each lane

Figure 3.13: Block diagram of 10-lane pipelined architecture.

consists of: A) Decoupling queues; B) Register file access; C) ALU execution; D) Write-back from the ALU back to the register file. The pipeline is five stages, with one clock cycle being given to every stage except the execute stage, which requires two cycles. The 128-deep instruction queue is implemented as a simple synthesized shift register, pre-loaded with instructions set with the scan chain, providing the ability to test different applications on our architecture.

An expanded view of the register file/ALU block is shown in Fig. 3.14. The ALU portion is split into two clock cycles, pipelining between the multiply and add functions of the multiply-accumulate to allow for higher max clock frequency, as the worst case path in the design determined from static timing analysis was found to be in the multiply-accumulate. A 32-bit operand register provides one of the inputs to the 32-bit adder of the multiply-accumulate.



Figure 3.14: Detailed block diagram of RF/ALU portion of pipeline.

### 3.3.1   Physical Implementation

The chip was fully digital, using a design flow produced and maintained by the OSU VLSI group. The full design was coded in Ver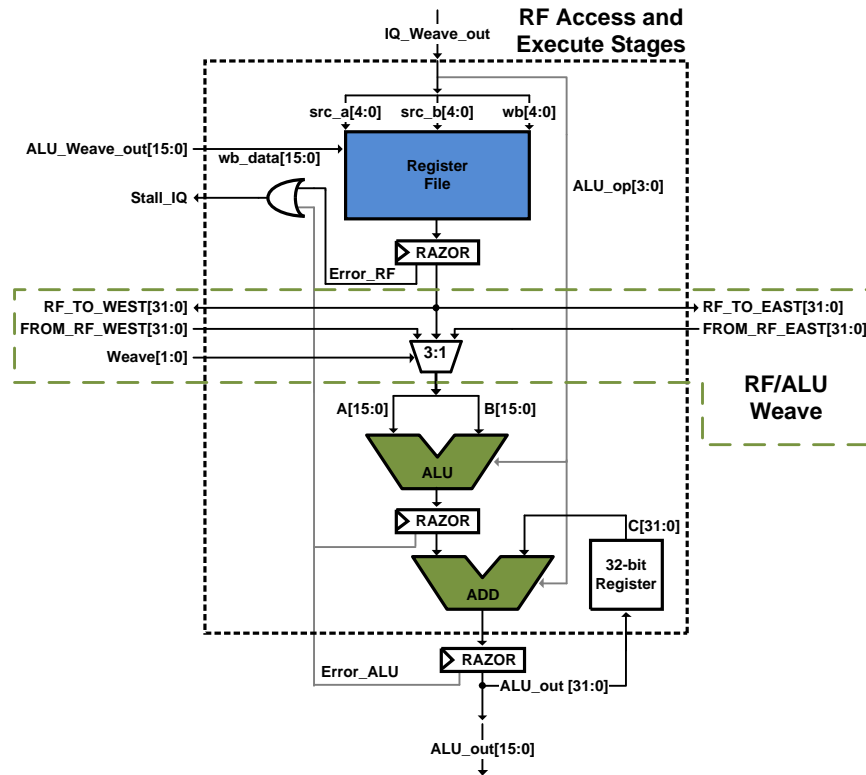ilog, and simulated for functionality using Modelsim. Synthesis was done using the Synopsys Design Compiler with an IBM 45nm SOI digital standard cell library characterized for a supply voltage of 0.9V. Place and route was performed with Cadence Encounter, and streamed into Cadence Virtuoso layout editor where the final DRC and LVS checks were done before being sent to the fab.

The top level layout for the Synctium test chip is shown in Fig. 3.15 and the fabricated chip is shown in Fig. 3.16. The The chip used 11 metal layers, and had 9 different power domains to allow for energy consumption measurements of all the separate blocks in the design. Each power domain can be identified in the image by the vertical yellow lines. The structure on the lanes is reflected in Fig. 3.16, showing that they are oriented horizontally across the supply domains. The power and ground input pads for each domain are paired across the top and bottom edges of the chip. Each component in the pipeline had its own scan inputs and outputs, to allow for characterization of each pipeline block. These scan input/output pads, clock pads, and enable signal pads are all place along the left and right edges of the chip, in a staggered two row formation. In total, there were 102 pads on the chip. The 2mm x 2mm chip was fabricated in a 45nm SOI process with a 1V nominal supply voltage, and functional from 1V down to 0.53V. Below 0.53V, the level shifters from the standard cell I/O library ceased to function properly.

Figure 3.15: Top level layout of Synctium test chip.

Extra attention was paid to the balancing of computational paths to ensure correct Razor error detection even under near-threshold operation. A negative consequence of this is the additional area and power overhead due to the insertion of timing buffers on the shorter delay paths. The area breakdown for the chip is shown in 3.1, showing that the Razor error detection circuits (including scanable pipeline flip-flops) was 27.52% of the total die area, and 34.58% of the area per lane. Traditional implementation of Razor (placing a detection circuit on

Figure 3.16: Die Photograph

only the worst case output paths) would have resulted in a more compact and power efficient design, but at the cost of higher susceptibility to unpredictably large timing delays in near-threshold. The area overhead of the decoupling queues and lane weaving circuitry is also an acknowledged drawback of the two resiliency methods. The decoupling queues exhibited an overhead of 13.43% per lane, and the lane weaving showed an overhead of 3.16%. It should be noted that these numbers are for a design that utilizes a small and simple ALU implementation. For a design with more-complex operations other than 16b fixed-point multiplies (i.e. double-precision floating point), ALUs would have taken up a much larger amount of area, making the overhead of the queues even less significant. Furthermore, the decoupling queues area overhead could be reduced significantly by using memory

compilers that result in more compact register files.

| Block | Instruction Queue | Decoupling Queues | Weaving | Register File | ALU | Razor + FF |
|---|---|---|---|---|---|---|
| **Total Area** | 20.4% | 10.7% | 2.5% | 16.3% | 22.6% | 27.5% |
| **Per Lane** | N/A | 13.4% | 3.2% | 20.5% | 28.3% | 34.6% |

Table 3.1: Area Breakdown

## 3.4 System Implementation

### 3.4.1 Board Design and Testing

The board designed for testing is shown in Fig. 3.17. A 121 pin PGA package was used, both purchased and wire bonded by Quik-Pak. A matching Zero-Insertion-Force socket was used to allow for measurements across many different chips. The PGA was chosen based on the requirement of having a sufficient number of pins while still remaining simple enough to design the board properly. While this PGA fit that requirement, the cavity of the package was much larger than the size of the chip, creating bond wire lengths of around 8mm. This could of potentially led to issues with the clock at higher frequencies, though because our target was for near-threshold operation at lower speeds, this did not serve as a limitation for us.

Also included on the test board were voltage dividers on all of the chip inputs to step the voltage down from the 3.3V output of the Ni-Daq measurement tool to the max 1.2V input to the pad-ring. Signal output from the chip get level shifted up to 3.3V before being input back into the Ni-Daq. Clock inputs (clock, delayed

Figure 3.17: Test board photograph

Razor clock, and clock for a separate test) were directly input via SMA connectors and laid out with thick traces to maintain as strong of a signal as possible at the chip.

Experiments were run using the following tools. A National Instruments Data Acquisition tool was used for scan-in and scan-out of data, clock signals were generated using a Textronix data timing generator, and all supply voltages were generated using 3, 3-output Agilent power supplies. Test code to control all of the equipment this unit was written in C-code using the National Instruments Labwindows/CVI test interface. The power supply and data timing generator were controlled through GPIB and the Ni-Daq was controlled through ethernet.

The test sequence was straightforward and fully automated. The goal of the testing was first to characterize each block's max frequency, then to run a set of

instructions pre-loaded into the instruction register, monitoring the error count and measuring the dynamic energy consumption. The characterization test involved first scanning in a known set of inputs to all blocks in parallel, clocking once, and then scanning out from the output end of each pipeline stage. If the output was correct, the clock speed was increased. Eventually, the top clock speed was found using a binary search. As previously mentioned, this calibration pattern could be time consuming, but the max time was cut down by testing all blocks in parallel. After all of the blocks were measured, the weave control bits were scanned in, along with a known set of instructions into the instruction register and clocked through the pipeline, monitoring the error rate and energy consumption.

## 3.5   Measurement Results

Figures 3.18 and 3.19 show both the measured delay distribution of a 16-bit multiply across ten lanes and the normalized delay variability across supply voltages, respectively. In Figure 3.18, each box-plot shows a 10-vector delay distribution, mean delay (black dot) with one standard deviation range (thick bar), and the delay of the worst-case path as determined by STA during synthesis, and verified with Hspice simulations on the post-layout extracted netlist. To justify running measuring only 10 vectors per lane at every supply voltage, 100 vectors were also measured at 0.53V. The smaller plot in the top right corner of Fig. 3.18 comparing the delay distributions for 10-vectors versus 100-vectors shows only a slight change in mean delay variation while maintaining same the best and worst-case

Figure 3.18: Delay Variation of all 10 lanes across multiple voltages.

delays. This shows that the 10-vector results will shift only slightly as the number of vectors tested increases.

A comparison against post-layout extracted Monte Carlo simulations shows that variation between lanes, especially at low voltages, is hard to predict (Fig. 3.19. It can also be observed that the delay spread is much greater at 0.53V versus 1.0V. At 0.53V the difference between the fastest and slowest lane is as much as 6ns, which results in about a 200% increase. At 1.0V the difference is just 0.4ns, resulting in a 20% increase from the fastest delay. The maximum frequency (limited by the STA worst case path) at 0.53V for a 16-bit multiply is 85MHz. It

should also be noted that the fast paths of all of the lanes remain well balanced (within 20% of each other), showing that the buffer insertion was effective, and the use of Razor, even at $V_{DD} = 0.53V$, is possible. At 1.0V, the delay distribution improves significantly, but the maximum throughput is limited to only 500MHz because of clock rise/fall time degradation on the PCB, bond wires, and package parasitics.



Figure 3.19: Mean normalized delay variability for all 10 lanes across multiple voltages

The performance improvement of the decoupling-queues enabled DPSP pipeline over a Razor-only SIMD pipeline is shown in Fig. 3.20. The plot shows the throughput trends of the two different pipelines as the error rate is increased. Denoted along the curves are two different sets of operations: a 16 bit matrix

multiply and a 256-point complex FFT. These operations were measured on chip, and compared with the curves, which were generated from a Matlab analytical model of expected performance. During testing, error rates were pre-determined using input vectors that were known to cause errors for a given application. The errors were then tracked using error counters that monitored each lane.



Figure 3.20: Throughput effect of decoupling queues.

The plot shows that depending on the error rate, up to a 19% throughput increase can be gained by using decoupling queues. Due to more frequent micro-barriers, applications such as FFT exhibit a slight decrease in performance, compared to an application such as matrix multiplication, when the error rate goes above 1%. At small error rates (less than 0.0001%), the throughput improvement

of DPSP over conventional SIMD is minimal. Conversely, when the error rate is larger than 12.5%, the throughput of both SIMD and DPSP saturates, as this is the crossover point where the likelihood of an error in any one instruction among the eight active lanes is 100%. The throughput difference between the two designs at the point of saturation is because the DPSP design only exhibits a one-cycle penalty for error propagation versus a five-cycle penalty in a Razor-only SIMD design using counter-flow pipelining recovery.



Figure 3.21: Lane weaving effect on frequency and throughput across multiple voltages.

Fig. 3.21 shows the measured throughput improvements with the use of lane weaving as supply voltage is decreased and timing variations increase. In order for

a throughput improvement to be observed, the increase in maximum operating frequency must be greater than the fraction of lanes eliminated. That is, the number of active lanes multiplied by the new increased freque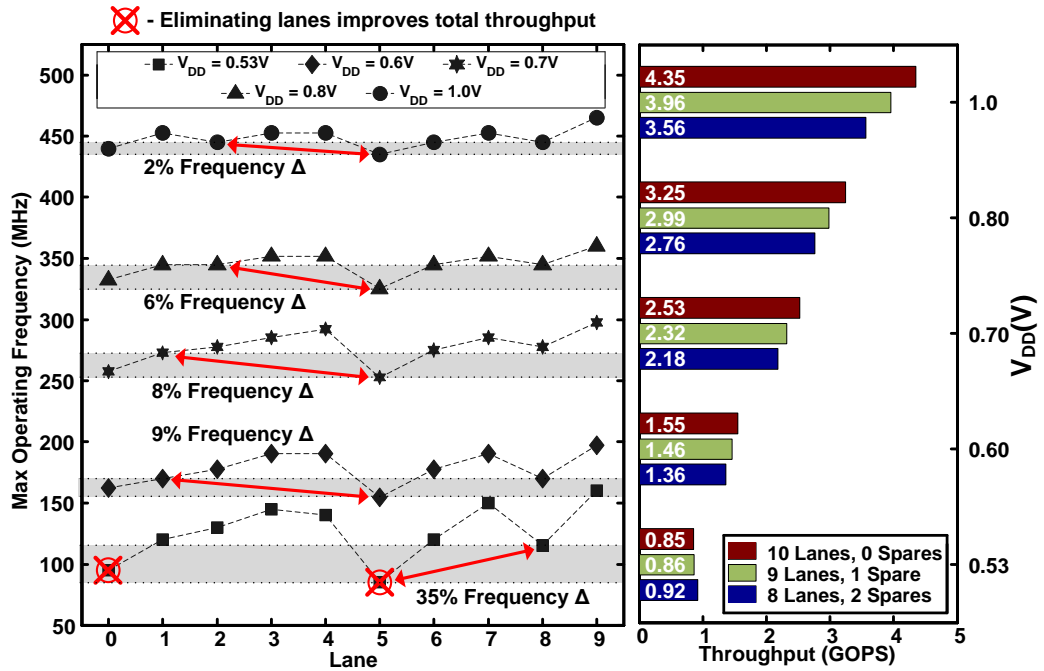ncy must be larger than the lower frequency multiplied by the maximum number of lanes. For $V_{DD}$ = 0.6-1.0V, no throughput improvement is observed because the maximum operating frequency improves by less than 10% when the components with the worst delay are disabled. However, at $V_{DD}$ = 0.53V, the elimination of the one or two worst components (i.e. disabling one or two lanes) enables an increase in operating frequency of 11% and 35%, respectively. This results in a throughput improvement of 1% and 8% for one-spare and two-spare configurations, respectively. Utilizing a two-spare setup with variation resiliency at $V_{DD}$ = 0.53V resulted in an energy/operation decrease of 3.5% over zero-spare and 2.6% over one-spare configurations.

The performance and energy-efficiency of the processor at $V_{DD}$=1.0V and $V_{DD}$=0.53V is summarized in Table 3.2. Measurement results are shown for a 16-bit matrix multiply application with a known error rate of 0.5%. Frequency and throughput values are taken for three different error resiliency cases: Without any resiliency techniques; with Razor on a SIMD pipeline; and with Razor on a DPSP. The results show that inclusion of only Razor circuits on a traditional SIMD pipeline with 5 cycle error recovery using counterflow pipelining will increase throughput by  0.5%, however the energy efficiency actually decreases by  25% at both 1.0V and 0.53V. Using Razor in combination with DPSP results in a much larger throughput gain of 27% over the baseline condition. Measurement results also show an increase in GOPS/mW over Razor-only designs of 9.8% at 1.0V and

11.2% at 0.53V, however, they still result in a decrease of 15% when compared to the pipeline with no error resiliency techniques.

The decrease in GOPS/mW of both the Razor-only and the Razor with DPSP design when compared with the baseline numbers can be explained by examining the energy/operation/lane breakdown in Table 3.2. The Razor logic accounted for a significant portion (24%) of the total energy consumption, this was due mostly to the need to include Razor flip-flops on every output bit of every stage, as well as adding a large number of buffers to balance the timing of the computational paths in the logic. While the throughput did increase by adding the Razor error

| Technology | 45nm SOI | |
|---|---|---|
| Architecture | 8-Lane SIMD with 2 spare lanes | |
| Die Area | 2mm x 2mm | |
| Application | 16-bit Matrix Multiply with Error Rate of 0.5% | |
| Supply Voltage | 1.0 V | 0.53 V |
| | | |
| | **Without Razor or Variation Resiliency** | |
| Frequency | 435MHz | 85MHz |
| Throughput | 3.48GOPS/181.1mW | 0.68GOPS/10.13mW |
| | **With Razor** | |
| Frequency | 530MHz | 101MHz |
| Throughput | 3.64GOPS/249.63mW | 0.712GOPS/13.27mW |
| | **With Razor and Variation Resiliency** | |
| Frequency | 550MHz | 144MHz |
| Throughput | 4.4GOPS/274.6mW | 1.152GOPS/19.3mW |
| | | |
| | **Energy/operation/lane** | |
| 16b-Multiply | 17.4 pJ | 4.81 pJ |
| Decoupling Queue | 11.4 pJ | 3.16 pJ |
| Register File | 12.4 pJ | 3.58 pJ |
| Razor Logic | 14.5 pJ | 3.98 pJ |
| Clock distribution/Lane | 406fJ | 124 fJ |
| Leakage Energy | 5.12 pJ | 844 fJ |
| Total Energy | 61.226 pJ | 16.5 pJ |

Table 3.2: Performance Summary

detection and decoupling queues to the design, the energy overhead was too great, and unfortunately resulted in a decrease in GOPS/mW.

### 3.5.1 Design Analysis

The chapter described the design and implementation of a 10-lane near-threshold SIMD architecture with two error resiliency methods: lane weaving for static variations, and the decoupled parallel SIMD pipeline to combat dynamic variations. The design was fabricated in an IBM 45nm SOI process, and tested at 1.0V and 0.53V supply voltages. At $V_{DD}$=0.53V, the chip operated at a clock frequency of 144MHz, achieving 59.7 GOPS/W for a 16-bit multiply operation. At $V_{DD}$=1.0V, the chip operated at 550MHz, and achieved 16.02 GOPS/W.

While showing that these error resiliency methods were effective for improving performance at near-threshold, they actually ended up having a negative effect on performance per watt, which is a more important metric for low-power designs. The reasons for this are quite obvious, and serve as an important obstacle for performing error resiliency at near-threshold. The first is that implementing the Razor error detection circuits in near-threshold required a Razor flip-flop to be placed on every output bit, due to increased timing variations. This, combined with added buffers to balance out the timing in all of the computational paths, resulted in much greater power overhead than expected. Additionally, the added power from the implementation of the decoupling queues/instruction sequencer and frequency increase made improving the performance per watt with the existing

implementation unsuccessful.

The most important improvement to make before error resiliency techniques are used in near-threshold is to find a more efficient detection method. Prior work shows few options that have been developed for near-threshold [63], [64]. Development and implementation of an effective near-threshold detection technique is essential for improving reliability through implementation of lane weaving and DPSP. Optimizing the number of spare components for pipeline weaving and the depth of the queues for DPSP are also important, as the area and power overhead of implementing these techniques were a bit high. Better optimization could be determined through more simulations, depending on the type of applications that are expected to be run on the processor. Taking extra design time to fulfill this task would decrease overheads, and increase the feasibility of implementing these techniques on a low-voltage processor.

# Chapter 4: Conclusion

As more applications emerge where low-power is an important goal, circuit designers are exploring the option of operating the supply voltage in the sub/near-threshold region. Before fully moving in this direction, the limitations of low-voltage design need to be completely understood. While much research has been focused on the decreased performance and increase sensitivity to PVT variations, the radiation sensitivity resulting from voltage scaling has not been well defined.

In chapter 2 a test platform for radiation characterization across supply voltages was introduced and implemented in a TSMC 65nm CMOS process. This test chip was intended to be the first in a series of chips implemented across processes, that will eventually provide valuable information to circuit designers about the tradeoffs between process, supply voltage, and soft error rate. This chip included commonly used 6T and 8T memory structures, as well as a variety of logic tests focusing on transient pulse propagation in NAND-based vs. NOR-based logic, the effect of inverter static noise margin on the occurrence of transient pulses, and pulse propagation distance vs. inverter size. Neutron experiments were performed at the Los Alamos Neutron Science Center, and alpha particle experiments were done at the Oregon State University radiation center. Memory neutron SER vs. $V_{DD}$ trends were observed for the 6T SRAM, showing a 6.45x increase in SER and

a 2.6x increase in multi-bit upsets when supply voltage scales from 1.0V to 0.33V. Combinational logic tests under alpha irradiation gave valuable information about SETs vs. SNM, inverter size, and standard cell type. Test results also showed that while logic SETs do not occur often in 65nm standard cells at nominal voltage, once $V_{DD}$ dropped below 0.5V, the SETs increased significantly, showing that the logic SER contribution at sub/near-threshold needs to be taken into consideration.

Moving beyond this understanding, the next issue is that of how designers will address errors due to these variations and radiation effects. If the additional performance/functionality loss from increased variation and SER is too great at low-voltages, the practicality of designs in this operation region disappears. Therefore, steps need to be taken to regain some of the throughput loss of low voltages by operating at slightly faster clock frequencies and detecting and correcting timing errors in combinational logic.

In chapter 3 a 10-lane near-threshold SIMD processor is described and implemented in an IBM 45nm SOI process. The processor utilizes two error resiliency methods: lane weaving for static variations, and decoupled parallel SIMD pipeline to combat dynamic variations. Razor timing error detection is combined with clock gate error recovery and the Decoupled Parallel SIMD Pipeline to achieve a throughput of 1.152 GOPS at 144MHz with a supply voltage of 0.53V. While the area and power overheads of this implementation were larger than originally hoped for, this processor demonstrates that it is possible to reduce power by operating at low-voltages, while still maintaining the necessary amount of throughput needed for many modern applications.

# Bibliography

[1] R. Dreslinski *et al.*, "Near-threshold computing: Reclaiming moore's law through energy efficient integrated circuits," in *Proceedings of the IEEE*, vol. 98, no. 2, Feb. 2010, pp. 253–266.

[2] D. J. et. al, "An implantable 64nw ecg-monitoring mixed-signal soc for arrhythmia diagnosis," in *International Solid-State Circuits Conference (ISSCC)*, February 2014.

[3] K.-T. T. et. al, "A 0.5v 1.27mw nose-on-a-chip for rapid diagnosis of ventilator-associated pneumonia," in *International Solid-State Circuits Conference (ISSCC)*, February 2014.

[4] T. L. et. al, "A 0.48v 0.57nj/pixel video-recording soc in 65nm cmos," in *International Solid-State Circuits Conference (ISSCC)*, February 2013.

[5] R. Naseer *et al.*, "Critical charge characterization for soft error rate modeling in 90nm sram," in *ISCAS*, May 2007, pp. 1879–1882.

[6] K. McKay and K. McAfee, "Electron multiplication in silicon and germanium," *Physical Review*, vol. 91, pp. 1079–1084, 1953.

[7] R. Baumann, "Radiation-induced soft errors in advanced semiconductor technologies," *IEEE Transactions on Device and Materials Reliability*, vol. 5, no. 3, pp. 305–316, 2005.

[8] S. Walstra and C. Dai, "Circuit-level modeling of soft errors in integrated circuits," *IEEE Transactions on Device and Materials Reliability*, vol. 5, no. 3, pp. 358–364, 2005.

[9] P. Dodd *et al.*, "Impact of technology trends on seu in cmos srams," *IEEE Transactions on Nuclear Science*, vol. 43, no. 6, pp. 2797–2804, December 1996.

[10] J. Fu, C. Axness, and H. Weaver, "Memory seu simulations using 2-d transport calculations," *IEEE Electronic Devices Letters*, vol. EDL-6, no. 8, pp. 422–424, 1985.

[11] P. R. et. al, "Determination of key parameters for seu occurrence using 3-d full cell sram simulations," *IEEE Transactions on Nuclear Science*, vol. 46, no. 6, pp. 1354–1362, December 1999.

[12] L. Freeman, "Critical charge calculations for a bipolar sram array," *IBM Journal of Research and Development*, vol. 40, no. 1, pp. 77–89, January 1996.

[13] T. Heijmen, D. Giot, and P. Roche, "Factors that impact the critical charge of memory elements," in *IOLTS*, July 2006, pp. 57–62.

[14] T. M. et. al, "Criterion for seu occurrence in sram deduced from circuit and device simulations in case of neutron-induced ser," *IEEE Transactions on Nuclear Science*, vol. 52, no. 4, pp. 1148–1155, August 2005.

[15] T. May and M. Woods, "A new physical mechanism for soft error in dynamic memories," in *International Reliability Physics Symposium*, 1978, pp. 33–40.

[16] J. Ziegler, M. Ziegler, and J. Biersack, "Srim - the stopping and range of ions in matter," *Nuclear Instruments and Methods in Physics Research*, pp. 1818–1823, 2010.

[17] B. Clark, M. Weiser, and I. Rasiah, "Alpha radiation sources in low alpha materials and implications for low alpha materials refinement," *Thin Solid Films*, vol. 462-463, pp. 384–386, 2004.

[18] D. Lambert *et al.*, "Neutron-induced seu in srams: Simulations with n-su and n-o interactions," *IEEE Transactions on Nuclear Science*, vol. 52, no. 6, pp. 2332–2339, 2005.

[19] F. Wrobel *et al.*, "Incidence of multi-particle events on soft error rates caused by n-si nuclear reactions," *IEEE Transactions on Nuclear Science*, vol. 47, pp. 2580–2585, 2000.

[20] J. Dirk *et al.*, "Terrestrial thermal neutrons," *IEEE Transactions on Nuclear Science*, vol. 50, no. 6, pp. 2060–2064, December 2003.

[21] J. Ziegler, "Terrestrical cosmic rays," *IBM Journal of Research and Development*, vol. 40, no. 1, pp. 19–39, 1996.

[22] R. Baumann, T. Hossain, S. Murata, and H. Kitagawa, "Boron compounds as a dominant source of alpha particles in semiconductor devices," in *International Reliability Physics Symposium*, 1995, pp. 297–302.

[23] S. W. et. al, "Thermal neutron soft error rate for srams in the 90nm-45nm technology range," in *International Reliability Physics Symposium*, May 2010, pp. 1036–1039.

[24] P. H. et. al, "Neutron soft error rate measurements in a 90-nm cmos process and scaling trends in sram from 0.25-um to 90-um," in *IDEM*, December 2003, pp. 21.5.1–21.5.4.

[25] H. Fuketa, M. Hashimoto, Y. Mitsuyama, and T. Onoye, "Alpha-particle-induced soft errors and multiple cell upsets in 65-nm 10t subthreshold sram," in *International Reliability Physics Symposium*, 2010, pp. 3A.4.1–3A.4.5.

[26] ——, "Neutron-induced soft errors and multiple cell upsets in 65-nm 10t subthreshold sram," *IEEE Transactions on Nuclear Science*, vol. 58, no. 1, pp. 2091–2102, August 2011.

[27] M. Casey *et al.*, "Single-event effects on ultra-low power cmos circuits," in *International Reliability Physics Symposium*, 2009, pp. 194–198.

[28] D. McMorrow *et al.*, "Subbandgap laser-induced single event effects: Carrier generation via two-photon absorption," *IEEE Transactions on Nuclear Science*, pp. 3002–3008, 2002.

[29] M. Gadlage *et al.*, "Digital device error rate trends in advanced cmos technologies," *IEEE Transactions on Nuclear Science*, vol. 53, no. 6, pp. 3466–3471, December 2006.

[30] B. Gill, N. Seifert, and V. Zia, "Comparison of alpha-particle and neutron-induced combinational and sequential logic error rates at the 32nm technology node," in *International Reliability Physics Symposium*, 2009, pp. 199–205.

[31] P. Jannaty *et al.*, "Two-dimensional markov chain analysis of radiation-indueced soft errors in subthreshold nanoscale cmos devices," *IEEE Transactions on Nuclear Science*, vol. 57, no. 6, pp. 3768–3774, December 2010.

[32] P. S. et. al, "Modeling the effect of technology trends on the soft error rate of combinational logic," in *IEEE DSN*, June 2002, pp. 389–398.

[33] L. Chang *et al.*, "An 8t-sram for variability tolerance and low-voltage operation in high-performance caches," *IEEE Journal of Solid-State Circuits*, vol. 43, no. 4, pp. 956–963, April 2008.

[34] M. Yamaoka *et al.*, "Low-power embedded sram modules with expanded margins for writing," in *IEEE Solid-State Circuits Conference*, 2005, pp. 480–481.

[35] B. Calhoun and A. Chandrakasan, "A 256-kb 65-nm sub-threshold sram design for ultra-low-voltage operation," *IEEE Journal of Solid-State Circuits*, vol. 42, pp. 680–688, March 2007.

[36] B. Narasimham, "Characterization of heavy-ion, neutron and alpha particle-induced single-event transient pulse width in advanced cmos technologies," Ph.D. dissertation, Vanderbilt University, December 2008.

[37] N. Seifert, *Radiation-induced soft error: A Chip-level Modeling Perspective.* now Publishers inc., 2010.

[38] T. Calin, M. Nicolaidis, and R. Velazco, "Upset hardened memory design for submicron cmos technology," *IEEE Transactions on Nuclear Science*, vol. 43, no. 6, pp. 2874–2878, December 1996.

[39] (2013). [Online]. Available: http://www.ti.com/product/lmk03806

[40] (2011). [Online]. Available: http://www.ti.com/product/cdclvp1208 (

[41] (2009). [Online]. Available: http://wnr.lanl.gov/newwnr/4FP30L-A/4FP30L-A.shtml

[42] *Genesys Board Reference Manual*, Digilent, 1300 Henley Ct. Pullman, WA 99163, May 2013.

[43] (2008). [Online]. Available: http://www.xilinx.com/tools/microblaze.htm

[44] *Measurement and Reporting of Alpha Particle and Terrestrial Cosmic Ray-Induced Soft Errors in Semiconductor Devices*, JEDEC Test Standard No. 89A, 2001.

[45] C.-K. Chen *et al.*, "A hardware-efficient vlsi implementation of a 4-channel ica processor for biomedical signal measurement," in *2011 IEEE International Conference on Consumer Electronics*, January 2011, pp. 607–608.

[46] T. Rintaluoma and O. Silvén, "Simd performance in software based mobile video coding," in *2010 International Conference on Embedded Computer Systems (SAMOS)*, July 2010, pp. 79–85.

[47] G. Gammie *et al.*, "A 28nm 0.6v low-power dsp for mobile applications," in *ISSCC Dig. Tech. Papers*, 2011, pp. 132–134.

[48] A. Wang and A. Chandrakasan, "A 180-mv subthreshold fft processor using a minimum energy design methodology," *IEEE Journal of Solid-State Circuits*, vol. 40, no. 1, pp. 310–319, 2005.

[49] S. Hanson *et al.*, "Exploring variability and performance in a sub-200 mv processor," *IEEE Journal of Solid-State Circuits*, vol. 43, no. 4, pp. 49–63, April 2008.

[50] ——, "Ultrlow-voltage, minimum-energy cmos," *IBM Journal of Research and Development*, vol. 50, no. 4/5, pp. 469–490, 2006.

[51] J. Tschanz *et al.*, "Tunable replica circuits and adaptive voltage-frequency techniques for dynamic voltage, temperature, and aging variation tolerance." in *2009 Symposium on VLSI Circuits*, Kyoto, Japan, June 2009, pp. 112–113.

[52] K. A. Bowman *et al.*, "Energy-efficient and metastability-immune resilient circuits for dynamic variation tolerance," *IEEE Journal of Solid-State Circuits*, vol. 44, no. 1, pp. 49–63, January 2009.

[53] D. Ernst *et al.*, "Razor: circuit-level correction of timing errors for low-power operation," *IEEE MICRO*, pp. 10–20, 2004.

[54] J. Crop *et al.*, "Error detection and recovery techniques for variation-aware cmos computing: A comprehensive review," *Journal of Low Power Electronics Applications*, vol. 1, no. 3, pp. 334–356, 2011.

[55] Y. Tamir and M. Tremblay, "High-performance fault-tolerant vlsi systems using micro rollback," *IEEE Transactions on Computers*, vol. 39, pp. 548–554, 1990.

[56] R. Krashinsky *et al.*, "The vector-thread architecture," in *ISCA*, 2004.

[57] M. Woh *et al.*, "Anysp: Anytime anywhere anyway signal processing," in *ISCA*, 2009.

[58] N. Drego, A. Chandrakasan, and D. Boning, "All-digital circuits for measurement of spatial variation in digital circuits," *IEEE Journal of Solid-State Circuits*, vol. 45, pp. 640–651, 2010.

[59] E. Krimer, R. Pawlowski, M. Erez, and P. Chiang, "Synctium: a near-threshold stream processor for energy-constrained parallel applications," *Computer Architecture Letters*, vol. 9, no. 1, pp. 21–24, January 2010.

[60] R. Pawlowski *et al.*, "A 530mv 10-lane simd processor with variation resiliency in 45nm soi," in *IEEE Solid-State Circuits Conference*, February 2012, pp. 492–494.

[61] I. Koren and A. Singh, "Fault tolerance in vlsi circuits," *Computer*, vol. 23, no. 7, pp. 73–83, 1990.

[62] S. Gupta *et al.*, "The stagenet fabric for constructing resilient multicore systems," in *IEEE MICRO*, 2008, pp. 141–151.

[63] E. Krimer, J. Crop, M. Erez, and P. Chiang, "Replication-free single-event transient (set) detection for eliminating silent data corruption in cmos logic," in *Silicon Errors in Logic - System Effects (SELSE)*, March 2013.

[64] J. Crop, R. Pawlowski, and P. Chiang, "Regaining throughput using completion detection for error-resilient, near-threshold logic," in *Design Automation Conference (DAC)*, June 2012.