

45-nm Radiation Hardened Cache Design

by

Jerin Xavier

A Thesis Presented in Partial Fulfillment
of the Requirements for the Degree
Master of Science

Approved August 2012 by
Graduate Supervisory Committee:

Lawrence Clark, Chair
Yu Cao
David Allee

ARIZONA STATE UNIVERSITY

December 2012

ABSTRACT

Circuits on smaller technology nodes become more vulnerable to radiation-induced upset. Since this is a major problem for electronic circuits used in space applications, designers have a variety of solutions in hand. Radiation-hardening-by-design (RHBD) is an approach, where electronic components are designed to work properly in certain radiation environments without the use of special fabrication processes. This work focuses on the cache design for a high performance microprocessor. The design tries to mitigate radiation effects like SEE, on a commercial foundry 45-nm SOI process. The design has been ported from a previously done cache design at the 90 nm process node.

The cache design is a 16 KB, 4-way set associative, write-through design that uses a no-write allocate policy. The cache has been tested to write and read at above 2 GHz at $V_{DD} = 0.9$ V. Interleaved layout, parity protection, dual redundancy, and checking circuits are used in the design to achieve radiation hardness. High speed is accomplished through the use of dynamic circuits and short wiring routes wherever possible. Gated clocks and optimized wire connections are used to reduce power. Structured methodology is used to build up the entire cache.

DEDICATION

Dedicated to *God*, my parents *Xavier Thomas* and *Margaret Xavier*

ACKNOWLEDGMENTS

First of all, I would like to thank God for his grace and divine interventions without which I would not have achieved this feat. I would also like to thank my parents and family for the support they have given throughout my tenure here. They have stood by me through the ups and downs of my life and I am really grateful for that.

I am deeply honored to have Dr. Lawrence Clark as my advisor and for the guidance and support that he has given me throughout my masters. An embodiment of smart ideas and techniques, working under him has made me a more knowledgeable and productive engineer. Whenever stuck with a problem, he would always have a fool proof solution and there has never been a time where he is out of ideas. I also thank Dr. David Allee and Dr. Yu Cao for being on my committee.

I thank Dan Patterson, Satendra Kumar Maurya and Srivatsan Chelleppa for their support in the research and for their technical and non-technical ideas without which the successful completion of the work would not have been possible. Also I would like to thank my friends Nishant Chandra, Vinay Chinti, Ashutosh Singrarur and Winnie Mathews for their support in both academic as well as non-academic ventures.

TABLE OF CONTENTS

	Page
LIST OF TABLES	viii
LIST OF FIGURES	ix
CHAPTER 1. INTRODUCTION	1
1.1 Radiation Environment in space	1
1.2 Radiation effects on circuits.....	4
1.2.1 Single Event effects (SEE) in CMOS	5
1.2.1.1 SEE Mechanism	6
1.2.1.2 Types of single event effects	9
1.2.2 Total ionizing dose.....	12
1.2.2.1 TID Mechanism.....	12
1.3 Radiation hardening	14
1.3.1 Shielding	15
1.3.2 Radiation hardening by process (RHBP).....	15
1.3.3 Radiation Hardening by Design (RHBD).....	16
1.3.3.1 Design techniques for mitigating SEE effects	16
1.3.3.2 Design techniques for mitigating TID effects	19
1.3.4 Mitigation of destructive SEE.....	20
1.4 Thesis Organization	21
CHAPTER 2. 90 NM CACHE DESIGN.....	22
2.1 Introduction to Cache.....	22

	Page
2.2 RHBD Cache Design Requirements, Assumptions and Approaches	28
2.3 Cache Design and Architecture.....	29
2.3.1 Data Array and Tag Array	31
2.4 “Staticizing” the 90nm Cache Design.....	35
2.4.1 The Dynamic way hit design	36
2.4.2 The Static way hit design.....	37
2.5 45nm RHBD Cache Design	41
2.6 IBM 45nm SOI issues.....	42
CHAPTER 3. CIRCUIT DESIGN.....	47
3.1 Cache Design Requirements	47
3.2 Data Array Design Details	47
3.2.1 SRAM Cell	49
3.2.2 Word Line (WL) Decoder.....	51
3.2.3 Write and Precharge Circuitry	53
3.2.4 Techniques to Achieve SEE Radiation Hardness	55
3.2.4.1 Error Checking Circuitry	59
3.2.4.1.1 WL Encoders and NOR Checker	60
3.2.4.1.2 Wren Encoder and NOR Checker	63
3.2.4.1.3 BL Precharge Checker	64
3.2.4.1.4 BL Read Checker	66
3.2.4.1.5 Write Checker	68

	Page
3.2.4.1.6 SRAM cells used in Read/Write Checkers	70
3.3 Tag Array Design Details	73
3.3.1 Tag Array Circuits	74
3.3.2 Hit Generation.....	76
3.3.3 Error Checking Circuits	76
3.3.4 Write and precharge circuitry	76
3.3.5 Dual Redundant Hit Generation	77
CHAPTER 4. ARRAY LAYOUT	79
4.1 Tag Bank Layout.....	79
4.1.1 Layout Interleaving.....	79
4.2 Data Bank Layout	80
4.2.1 Interleaved Layout and Parity Protection	81
4.3 Layout Techniques	81
4.4 Techniques to Reduce Power.....	82
CHAPTER 5. FULL CACHE LAYOUT	84
5.1 The structured methodology flow	84
5.1.1 Steps involved in the structured flow	87
5.1.2 Characterizing macro blocks	87
5.2 Cache layout.....	88
5.2.1 The data array layout	88
5.2.2 The tag array layout	90

	Page
5.2.3 Full cache layout	91
CHAPTER 6. CONCLUSION.....	95
REFERENCES	98

LIST OF TABLES

Table	Page
2.1 MAJOR INTEL [®] MICROPROCESSORS AND THE CACHE SIZE.....	27
3.1 AREA COMPARISON BETWEEN THE SRAM CELLS FOR 45 NM PROCESS.....	51
3.2 THE TRUTH TABLE OF WL ENCODER PART A AND PART B.....	63
3.3 TRANSISTOR SIZES USED FOR THE LOGIC SRAM CELL.....	71
3.4 READ MARGIN FOR THE DUMMY SRAM CELL FOR 45 NM PROCESS.....	71
5.1 DELAY AND AREA FOR VARIOUS CACHE DESIGNS	93
5.2 AVERAGE POWER	94
5.3 ENERGY PER OPERATION	94

LIST OF FIGURES

Figure	Page
1.1 Cartoon showing the space radiation environment. (After [58] [59])	4
1.2 (a) Ion strike at the output of an inverter. (b) Funnel formation and charge	7
1.3 Funneling in an n+/p silicon junction following an ion strike showing	8
1.4 Band diagrams showing the process of TID damage. After([57]).....	14
1.5 TID induced leakage in MOS devices. (After [26])	15
1.6 Implementation of TMR based hardware scheme	17
1.7 Layout of an edgeless transistor with a P+ guard ring for inter-device TID induced leakage mitigation (After [37])	19
1.8 Layout of (a) edge transistor (b) edgeless transistor used in TID mitigation ..	19
2.1 Memory Hierarchy	22
2.2 A high level example of the Cache structure [47].	25
2.3 Miss rate for vs cache size and set associativity. (After [50])	26
2.4 The basic diagram of cache. (After[47]).....	30
2.5 Floor plan of cache. (After [47]).....	32
2.6 Basic diagram of data array (after [47]).....	33
2.7 Basic diagram of tag (after [47]).....	34
2.8 Overview of the Tag critical path. (After [47]).....	36
2.9 Hit Generation Circuits. (After [47])	37
2.10 Static way hit schematic.....	37
2.11 2x2 OAI Comparator static logic	38

Figure	Page
2.12 Simulation waveform for the cache critical path	39
2.13 Schematic of an Inverter SDL.....	40
2.14 Simulation snapshot for an inverter SDL.....	40
2.15 Modeling the parasitic neck loading due to the annular devices (after [47]).	41
2.16 Layout of the NMOS-access RHBD SRAM cell (after [47])......	41
2.17 SOI Inverter cross-section (after [52])......	42
2.18 IBM SOI process electron microscope snap (courtesy of IBM)......	43
2.19 Charge paths in an SOI device (after [52])......	44
2.20 Parasitic bipolar transistor in an SOI device (after [52])......	45
2.21 Pass gate leakage in dynamic gates (after [52])......	45
3.1 Basic diagram of data array.	48
3.2 Schematic of the SRAM cell.....	50
3.3 Two word lines are connected alternately to cells in one row.....	50
3.4 WL decoder.....	52
3.5 WL decoder simulation results.	53
3.6 Write and precharge circuitry. (After [47]).....	54
3.7 Self-timed circuits are used in WREN and precharge drivers.	54
3.8 Simulation waveform of WL, WREN and precharge.....	55
3.9 The error detection scheme of data array. (After [47]).....	57
3.10 The basic dynamic error checking circuit (after [42]).	60
3.11 The WL encode scheme [47].	61

Figure	Page
3.12 The WL encoder structure [47].	61
3.13 WL encoder and NOR checker (after [47]).	62
3.14 The WREN encoder and NOR checker (after [47]).	64
3.15 BL precharge suppression detection.	65
3.16 Simulation of an error detected on BL precharge.	66
3.17 Bit line read checker (after [47]).	67
3.18 Simulation of an error detected on BL during a read.	68
3.19 Write checker. (After [47])	69
3.20 Simulation of an error detected during a write.	70
3.21 Layout of the Logic SRAM cell.	71
3.22 Monte Carlo simulations for worst case read SNM analysis	72
3.23 Worst case write margin of the Logic SRAM cell.	73
3.24 Basic diagram of tag.	74
3.25 Schematic of the tag critical path.	75
3.26 Schematic of the precharge and wren circuitry.	77
3.27 Dual match generation of tag (after [47]).	77
4.1 Snapshot of the tag bank of the tag array.	79
4.2 Snapshot of the data bank of the data array	81
4.3 Gated clock in the generation of sub-bank clock (after [47]).	83
5.1 Structured methodology flow (after [53]).	85
5.2 Data banks and standard cell placement	89

Figure	Page
5.3 Fully routed data array	89
5.4 Tag banks and standard cell placement	91
5.5 Fully routed tag array	91
5.6 Floorplan of the cache.....	92
5.7 The routed entire Cache	92
5.8 HSPICE simulation showing the Cache critical path	93

Chapter 1. INTRODUCTION

Electronic systems operating in radiation environment are increasingly vulnerable to radiation effects, due to fabrication process scaling, decreasing feature sizes, supply voltages and lower noise margins. Single event effects (SEEs) are caused when radiation particles such as protons, neutrons, alpha particles, or heavy ions strike sensitive diffusion regions of transistors in circuits. Studies have shown that SEEs are troublesome for military and space applications and radiation-induced single-event transients (SET) were the primary failure mechanism behind several spacecraft malfunctions in recent years [1-4]. There are also critical applications like biomedical, industrial and banking systems that also demand highly reliable systems [5]. Consequently, in recent years, the study and analysis of radiation effects on circuits has been a major research area. The technique of designing and fabricating electronic systems to withstand radiation is called radiation hardening [6]. This chapter provides an overview of the radiation environment, radiation effects on electronic devices and circuits and relevant techniques to achieve radiation hardness.

1.1 Radiation Environment in space

The space environment harbors phenomena that can be potentially hazardous to human and technological systems. This environment is not static and includes variations caused by solar flares or coronal mass ejections. This combined environment creates a multitude of issues for electronic systems contained within space systems [7]. The complex radiation spectrum of the space

environment typically consists of charged particles originating from various sources.

Solar energetic particles (SEP) are large fluxes of atomic particles, primarily electrons and protons with energies of the order of MeV that are accelerated and expelled from the sun by its solar flares. Solar flares generally vary on the scale of minutes to a few days, in response to events such as storms or sub-storms in the Sun. Trapped particles, which are 93% protons, 6% alpha particles, and about 1% heavy nuclei, contribute the most to radiation effects in low and medium Earth orbits that pass through the Van Allen belts [8]. The Van Allen belt(s) comprises a ring of particles trapped by the earth's magnetic field and consist of mostly high energy (1-10MeV/nucleon) electrons [9]. Figure 1.1 illustrates the space environment and the Van Allen belts with respect to the earth. Galactic cosmic rays (GCR) comprise of ions from almost all elements in the periodic table. Even though typically found in much lower fluxes than trapped particles, they can have energies as high as a TeV/nucleon. These GCRs are modulated by the 11-year solar cycle or sunspot cycle [7]. GCRs are about 87% protons, 12% helium, with the remainder composed of heavy ions through actinides [11]. CRAN particles (primarily cosmic ray albedo-neutrons) are primarily secondary cosmic ray neutrons produced by the interaction of GCR with the earth's atmosphere at about 55km above the earth surface. These have a half-life of 11.7 minutes beyond which they decay in to an electron, proton and an anti-neutrino. Secondary neutrons produced by the interaction of interplanetary

GCRs and solar particles with the atmosphere, are the most important contributor to single event effects at altitudes below 60,000 feet. The rest of the electromagnetic spectrum in space consists of X-rays (wavelengths $10\text{\AA} - 100\text{\AA}$), extreme ultraviolet or EUV ($100\text{\AA} - 1000\text{\AA}$), ultraviolet ($1000\text{\AA} - 3500\text{\AA}$), the visible spectrum ($3500\text{\AA} - 7000\text{\AA}$) and the infra-red spectrum ($0.7\mu - 7\text{mm}$). Each type of radiation has a characteristic spectrum and preferred interaction mode with matter that give rise to various effects such as photo-ionization, photoelectron emission, Compton effect, etc. Photon interactions are not a primary concern for satellites in the natural space environment [11]. Only those phenomena from the vast and complex radiation effects in space that affect circuits would be the focus of this research.

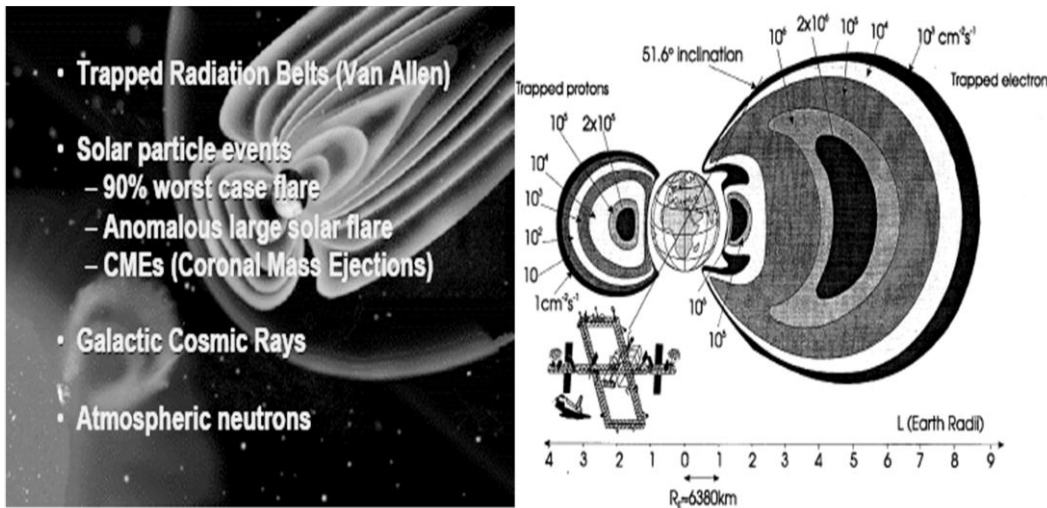


Figure 1.1 Cartoon showing the space radiation environment. (After [58] [59])

1.2 Radiation effects on circuits

Radiation effects can lead to degradation, malfunction, or even permanent damage in electronic circuits and devices [12]. Radiation interactions with solid material depends on a variety of factors like the material type, kinetic energy, mass, and charge state of the incoming particle and the mass, atomic number and density of the target material to name a few.

An ion travelling through a material loses its kinetic energy primarily through coulombic interactions with the electrons of that material leaving a trail of charge in its path. The higher the energy of the ion, the longer the distance it travels before being stopped by the material. The stopping power or linear energy transfer (LET) is a function of the material through which a charged particle is traveling and refers to the energy loss of the particle per unit length in the

material. The LET (MeV-cm²/mg) is a function of both the ion's mass and energy and density of the target material.

$$LET = \frac{1}{\rho} \frac{dE}{dx} \quad (MeV - cm^2/mg) \quad (1)$$

where $\frac{dE}{dx}$ is the energy loss per unit length and ρ is the material density in mg/cm³. The maximum LET value near the end of the particle's range is called the Bragg peak [13]. Bragg peak is the peak on the Bragg curve which plots the energy loss of ionizing radiation when it travels through matter.

Radiation particles interact with material, depositing charge by two major mechanisms: direct ionization and indirect ionization. In direct ionization, a high energy charged particle interacts directly with the electrons in the target material, breaking them free from their bound states, creating a dense track of free charge. During indirect ionization, the high energy particle collides with a nucleus in the target material, freeing that nucleus from its bound location. This recoiling nucleus is the charged particle which then creates the charge track. The two major radiation effects on MOS circuits and devices are single event effects (SEEs) [14] and total ionizing dose (TID) effects [15].

1.2.1 Single Event effects (SEE) in CMOS

SEEs are caused by a single radiation particle strike. All single-event effects are caused by the collection of charge at a sensitive region of an electronic device following the passage of an energetic particle through the device as shown

in Figure 1.2(a). Heavy ions, protons and neutrons constitute the majority of the particles responsible for this effect. Radiation effects from heavy ions are most often due to direct ionization while the vast majority of SEEs from protons are due to indirect ionization through collisions with heavier nuclei. Except for certain devices like photo detectors that are designed to detect small amounts of charge and large collection lengths, direct proton ionization is a rare phenomenon in space. This is due to the fact that protons have small LET owing to their smaller mass. For example 60 M-eV protons have an LET as low as $0.008 \text{ M-eV.cm}^2 / \text{mg}$ in Si for a collection depth of $10\mu\text{m}$ [56]. SEEs from neutrons are entirely due to indirect ionization, as they do not cause direct ionization owing to their neutral charge [16]. This is due to the fact that protons have small LET owing to their smaller mass. For example 60 M-eV protons have an LET as low as $0.008 \text{ M-eV.cm}^2 / \text{mg}$ in Si for a collection depth of $10\mu\text{m}$ [56]. SEEs from neutrons are entirely due to indirect ionization, as they do not cause direct ionization owing to their neutral charge [16].

1.2.1.1 SEE Mechanism

Formation of a SEE involves three stages – charge generation, charge collection and circuit response. Charge generation is influenced by the particle's mass and energy and the properties of the materials it passes through. Charge is generated from a single event phenomenon generally within a few microns of the junction. In silicon one electron-hole pair is produced for every 3.6 eV of energy

lost by the incident radiation. As silicon has a density of 2.328 g/cm^3 , it is easy to calculate from equation (1) that an LET of $97 \text{ MeV-cm}^2/\text{mg}$ corresponds to a charge deposition of 1 pC/sq.m . Hence the amount of collected charge in silicon can be given by the formula

$$Q = 0.01036LET \quad \text{pC} / \mu\text{m} \quad (2)$$

Thus, the charge collected (Q) for these events range from one to many hundreds of fC depending on the type of ion, its trajectory, and its energy over the path through or near the junction. The worst-case would be in case the junction is floating (as in DRAMs, dynamic logic circuits without keepers, and some analog designs) and is extremely sensitive to any charge collected from a radiation event. As discussed above when a particle strikes a microelectronic device, the sensitive regions are the reverse biased p/n junctions, as illustrated in Figure 1.2(b). Charge generated along the particle track can locally collapse the junction electric field

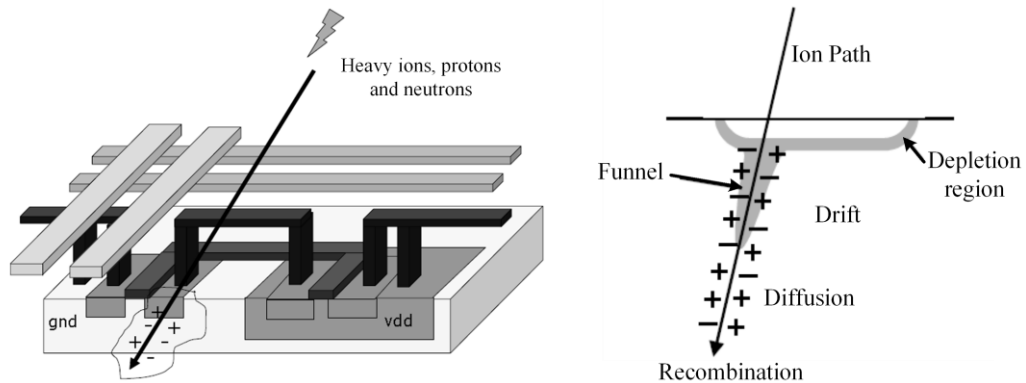


Figure 1.2 (a) Ion strike at the output of an inverter. (b) Funnel formation and charge collections mechanisms in the semiconductor following an ion strike. (After [18])

due to the highly conductive nature of the charge track creating a “field funnel” as shown in Figure 1.3 [18]. This funneling effect increases charge collection at the affected node by extending the junction electric field away from the junction deep into the substrate, such that charge deposited some distance from the junction gets collected through the drift process. This gives rise to a transient current at the junction contact. Strikes near a depletion region can thus result in a significant current transient as carriers diffuse into the vicinity of the depletion region field where they can be efficiently collected.

Even for direct strikes, diffusion plays a major role as carriers generated beyond the depletion region can diffuse back toward the junction. This can lead to unwanted current flow in associated circuit nodes.

Funneling is just not dependent on a direct strike on a depletion region. Near misses can also cause funneling if a high enough carrier density diffuses into

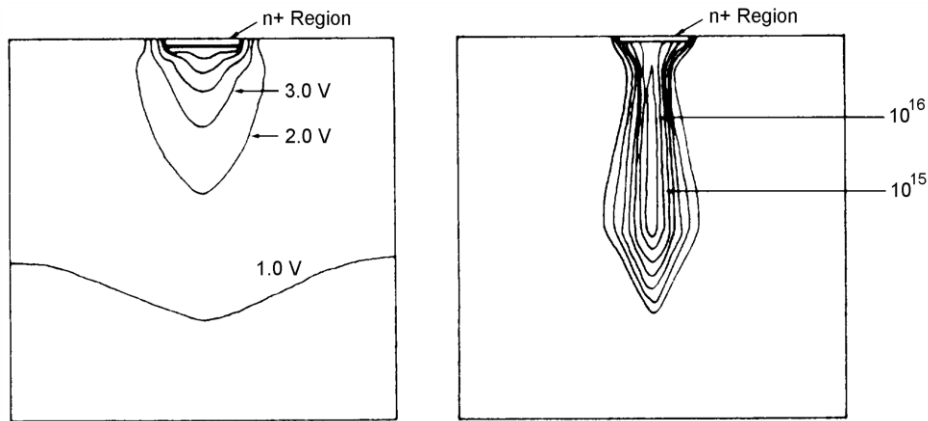


Figure 1.3 Funneling in an n+/p silicon junction following an ion strike showing contours of (a) electrostatic potential, (b) electron concentration (After [17]).

the depletion region to collapse it [18]. Due to the differences in the hole and electron mobility, funneling occurs in reverse biased n/p diodes, but is much weaker or nonexistent in equivalent p/n diodes as shown in Figure 1.3. The applied voltage at the struck junction is not a constant and in fact the struck node may switch from being reverse-biased to zero-biased. In such cases, funneling may play a role in the early-time response of the circuit by initially helping to flip the node voltage, but it is the late-time collection by diffusion that ensures the node stays flipped. This causes an error in the circuit operation.

The device characteristic that determines the sensitivity of a device to radiation upsets is its critical charge. This is the amount of charge needed at the terminal of the device to flip the node state.

1.2.1.2 Types of single event effects

Single-event effects can be characterized as non-destructive (causing a soft error) [19] or destructive SEE (resulting in a hard error). The error is called “soft” because the circuit is not permanently damaged by the radiation – i.e. if new data is written to the bit cell, the device will store it correctly. In contrast, a “hard” error leads to the device being physically damaged such that it malfunctions, data is lost, and the damaged state is permanent.

Examples of non-destructive SEE include single-event transients, single event upsets in memory circuits (SEU), multi-bit upsets (MBU) and single event functional interrupts (SEFI). Destructive SEE include such phenomena as single-event latch up (SEL, which can be either destructive or non-destructive depending

on circuit design), single-event burnout (SEB), and single-event gate rupture (SEGR).

A single event transient (SET) is defined as a momentary voltage spike at a node in an integrated circuit [20-22]. The voltage spike can propagate through the combinational logic away from where it was generated and eventually appear at the circuit's output, given certain conditions. It may also be captured locally if it is generated within a latch, or non-locally if it first propagates through the circuit before being captured by a latch. Once the SET is captured by a latch or flip-flop it becomes a single event upset (SEU), and it is impossible to distinguish SEUs that result from SETs that have propagated from other locations in the circuit from SEUs that have been generated within the latch or flip-flop itself [23].

When a charged particle strikes one of the sensitive nodes of a memory cell, such as a drain in an off state transistor, it generates a transient current pulse that can turn on the gate of the opposite transistor. This effect can cause a bit flip in the memory cell causing a SEU in the memory cell [16]

If the radiation event generates sufficient charge, more than one bit can get affected, causing a multi-bit upset (MBU). MBUs are defined as the occurrence of two or more bit upsets, appearing within the same clock cycle from a single particle hit, to distinguish from random multiple hits within a single cycle [24]. While MBU usually constitute a fraction of the total observed SEU rate, their occurrences have significant implications for memory architecture in systems utilizing error correction methods.

Another type of soft error occurs when the bit is flipped in system control registers, such as that found in field programmable gate arrays (FPGAs) or DRAM control circuitry, so that the error causes the circuit to malfunction. This type of soft error, called a single event interrupt (SEFI) [25], impacts the product reliability since each SEFI leads to a direct product malfunction as opposed to typical memory soft errors that may or may not affect the final product operation depending on the algorithm, data sensitivity, etc. Digital functions most likely to be affected by SEFIs are clock and trees, PLLs, counters, address and control registers as well as poorly regulated power networks.

Single event latch-up is caused by a steady high current state that results when a parasitic silicon controlled rectifier (SCR) (p-n-p-n) structure is triggered into regenerative forward bias [17]. If the latch-up current is large enough this can lead to a destructive event.

Single event gate rupture (SEGR) and the single event breakdown (SEB) [17] are both mechanisms that are destructive and lead to hard failures. In SEGR the gate oxide is in a high conduction state (breakdown) initiated by a hit to the gate region while in SEB the junction is broken-down when the event causes avalanche and thermal runaway.

The rate at which soft errors occur is called soft error rate (SER). The unit of measure associated with SER is failure in time (FIT). One FIT is equivalent to one failure in 10^9 device hours.

1.2.2 Total ionizing dose

Total ionizing dose effects in electronics are caused by the interaction of radiation and the silicon dioxide insulating layers of the device [15]. When an MOS transistor is exposed to high-energy ionizing radiation, electron-hole pairs are created uniformly throughout the oxide. For MOS device degradation, the primary concern is electron-hole pair (ehp) generation in oxides (SiO_2) which leads to almost all total dose effects. The generated carriers induce charge build-up in the oxide and interface traps, which causes a threshold voltage shift and can lead to device degradation.

The total amount of energy deposited by particles that result in ehp production is commonly referred to as total ionizing dose (TID). The typical unit of TID that is used is rad, which denotes the energy absorbed per unit mass of SiO_2 .

1.2.2.1 TID Mechanism

Ionization in a target material is caused by the interaction of protons, electrons, energetic heavy ions, and photons with the atoms of that material. The unit used to measure TID is rad.

Radiation-induced charging of oxide involves many different physical mechanisms, which take place on very different time periods under different physical conditions. The physical processes (shown in Figure 1.4) that play a major role from the initial deposition of energy by ionizing radiation to the creation of ionization defects are:

1) Generation of ehp: Radiation-induced ionization damage is primarily the result of the generation of electron-hole-pairs (ehp) along the track of high energy electrons generated as a result of interactions by photons and protons. The density of e-hp generated along the tracks of incident particles is proportional to the LET of the incident particle and the band gap of the target material.

2) Prompt recombination of a fraction of the generated e-hp: Once generated, a fraction of the ehp are annihilated through recombination. The holes, which escape initial recombination, are relatively immobile and remain near their point of generation, where they form fixed charge in the oxide and cause a negative threshold voltage shift. The electrons generated in the oxide are relatively mobile and drift out of the oxide [26].

3) Transport of free carriers remaining in the oxide: This transport of holes to the Si/SiO₂ interface causes a short-term recovery of the threshold voltage.

and either

4a) Formation of trapped charge: As holes approach the Si/SiO₂ interface, some fraction of the holes get trapped in the device defects, forming a positive oxide-trap charge.

or 4b) formation of interface traps via reactions involving hydrogen in the SiO₂: Hydrogen ions (protons) are likely to be released as holes, “hop” through the oxide or get trapped near the Si/SiO₂ interface and form interface traps.

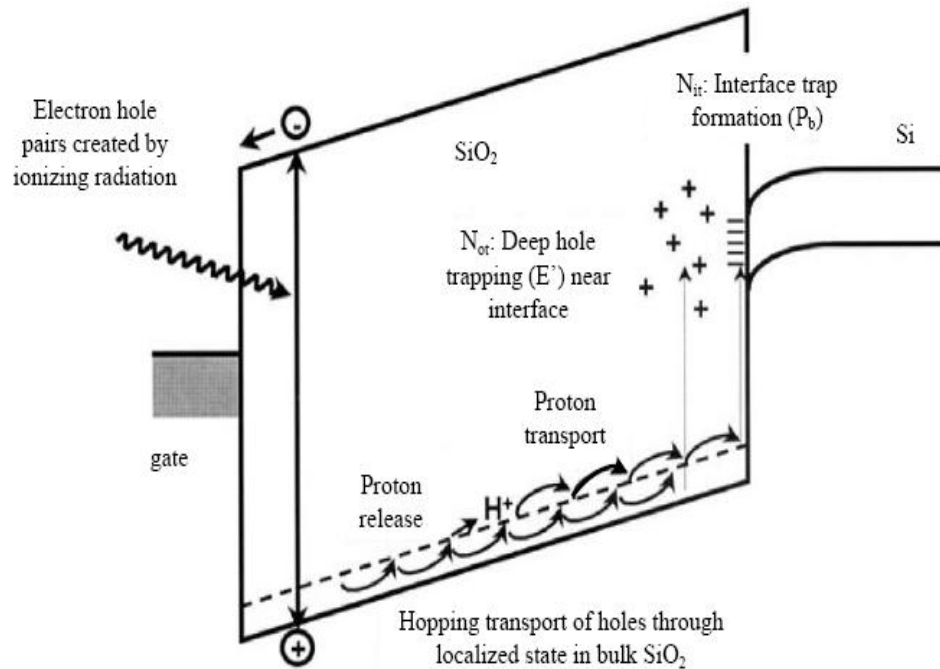


Figure 1.4 Band diagrams showing the process of TID damage. After([57])

This entrapment of holes in the oxide can shift threshold voltages of NMOS/PMOS transistors. In addition to oxide-trapped charge and interface-trap charge buildup in gate oxides, charge buildup will also occur in field oxides as shown in Figure 1.5.

1.3 Radiation hardening

There are several methods used for mitigating the effects of radiation. Radiation effects can be mitigated by using design techniques at all levels of the system design [12]. From the device level to the circuit level to the system level, there are methods that can be implemented to mitigate all types of radiation and all types of radiation effects.

There are many published techniques in soft error mitigation; the most common techniques are outlined in the subsequent sections.

1.3.1 Shielding

One simple method to mitigate the effects of radiation is the use of shielding. Since a majority of cosmic rays have very high energies, shielding has very little effect on them. For this reason shielding is seldom considered for the mitigation of SEE in electronic circuits.

1.3.2 Radiation hardening by process (RHBP)

Radiation hardening by process (RHBP) is a method to harden a device to TID and/or SEE using certain features in the fabrication process. This is done by modifying a current fabrication process. Modification is typically done by adding or changing process steps, ideally without impacting the performance or normal operating characteristics of the device. However higher costs makes it a less

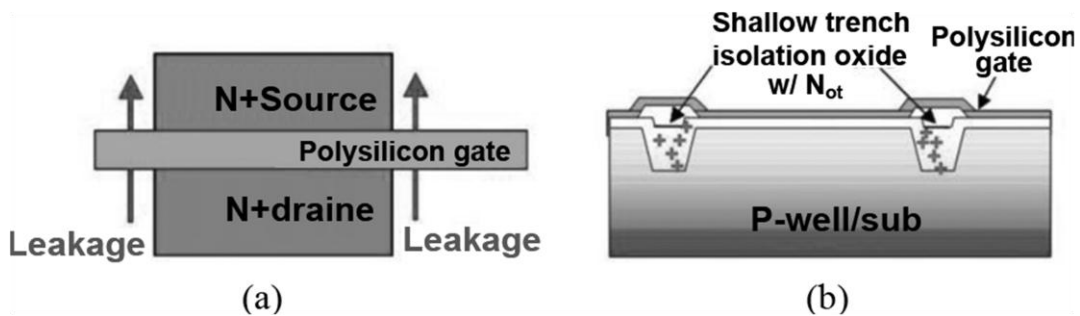


Figure 1.5 TID induced leakage in MOS devices. (After [26])

attractive approach.

1.3.3 Radiation Hardening by Design (RHBD)

Radiation hardened by design (RHBD), has recently become possibly the most popular approach uses design techniques implemented in a standard commercial foundry to make a non-hardened process hard to some level. RHBD techniques promise to improve the performance of rad-hard ICs by utilizing state-of-the-art commercial foundry silicon processes [27]. The work carried out in this thesis and presented in this report is based on these RHBD principles.

1.3.3.1 Design techniques for mitigating SEE effects

- Layout and Electrical level based techniques:
 - Built-in sensors for ionization detection: The bulk-built in current sensor (BICS) works as a monitor that senses the current at the bulk terminal. During fault-free operation, the current in the bulk is approximately zero. When a charged particle generates a current in the bulk, it is sensed by the BICS and the system control logic is notified to perform some fault tolerant technique to mitigate the effect.
 - Transistor resizing for charge dissipation: Widening transistors increases the capacitance of the most sensitive nodes making it harder for an SET to upset it [28].
- Logic-level based techniques:
 - Hardware (Spatial) redundancy using majority voting: A tri-voting system compares the outputs of three identical devices bit by bit, relying on the fact that while each bit is equally vulnerable to upset, the

probability of the same bit upsetting in two independent devices is very low [29] [30]. To save area and power, duplex architectures [54] has been developed. Two sets of circuits are calculated independently. Each of them has a parity bit to detect a single error. If there is an error in one set of circuit, the result of the other circuit is used. If both circuits have an error, a calculation error will be reported.

- Time redundancy using temporal filtering: Temporal voting is employed against SETs in which the same device or data path is polled 3 times and the results stored and voted as shown in Figure 1.6. One of the side effects of using this technique is that it limits the maximum speed at which the circuit can operate. If the delay is long enough it will literally filter out the signal [14] [31].

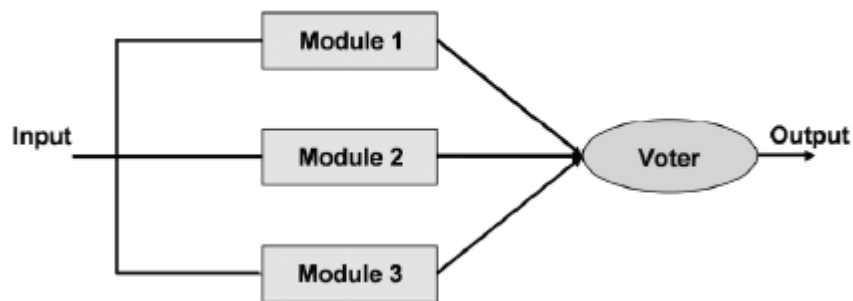


Figure 1.6 Implementation of TMR based hardware scheme

- Error correcting codes (ECC) for detection and correction of bit-flips [32]: In an error detection and correction (EDAC) scheme, redundant

bits are added to a data word to enable the system to detect and correct errors in the data (caused by SEU or SEFI) using ECC schemes such as Hamming codes [33]. Owing to the delay introduced by conventional EDAC, only large memories and L2 caches have used it in high performance ICs [54].

- Hardened memory cell to avoid bit-flipping: Memory elements can be protected against SEU by modifying their original design by including extra resistors or transistors. These extra transistors would be able to recover the stored value if a particle strikes one of the drains of a transistor in “off” state [34] [35].
- System level techniques:
 - Recovery and recomputation: Some highly reliable microprocessor systems maintain checkpoints that detect faults at various stages and try to recover the information [36]. Forward error recovery is detecting an error and continuing on in time, while attempting to mitigate the effects of the faults that may have caused the error. Backward error recovery is detecting an error and retracting back to an earlier valid system state or time. The later is more commonly used.

1.3.3.2 Design techniques for mitigating TID effects

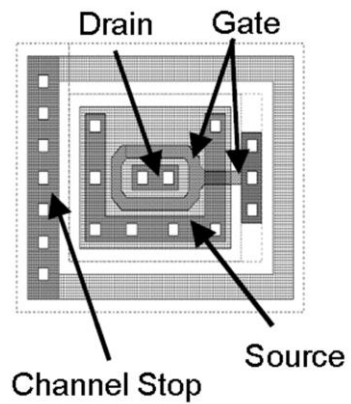


Figure 1.7 Layout of an edgeless transistor with a P+ guard ring for inter-device TID induced leakage mitigation (After [37])

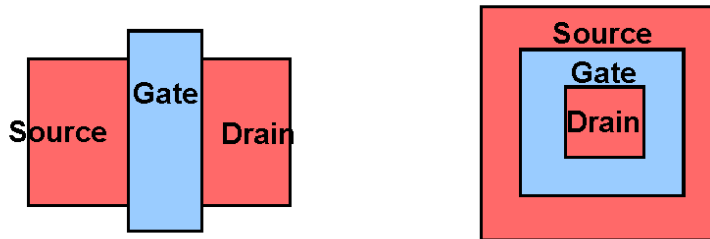


Figure 1.8 Layout of (a) edge transistor (b) edgeless transistor used in TID mitigation

The primary effect of total-dose irradiation on isolated, standard non-hardened layout NMOS transistors is to cause an increase in the off-state leakage

current. This increased leakage current is caused by the inversion of parasitic transistors at the transistor edges at or near the gate-oxide/field-oxide interface. A layout design technique that eliminates edge leakage by removing the parasitic edge transistors between the source and the drain is shown in Figure 1.8 [37]. This layout has no active diffusion edges overlapped by polysilicon that separates the source and the drain. This type of transistor with a closed-geometry layout is called an edgeless or annular transistor. TID induced inter-device leakage is caused when the field-oxide inverts due to exposure to ionizing radiation. The most common design solution to this problem is to surround each transistor with a p+ diffusion ring as shown in Figure 1.7.

1.3.4 Mitigation of destructive SEE

In contrast to SEU and SET mitigation, destructive single event effects would be hard to recover and it depends majorly on the individual device's response to the radiation effect. Hardening from the system level is difficult and in most cases, not very effective [38]. Non-recoverable destructive single event effects such as single event gate rupture and burn-out heavily damage devices in a manner which completely compromises the circuit's operation.

Detecting a hard error by using correction techniques like EDAC and by making use of unaffected bits to replace the damaged ones are effective techniques in mitigating destructive SEE effects.

1.4 Thesis Organization

This chapter provided a brief overview of the radiation environment, the effects of radiation on circuits and common radiation hardening techniques. Chapter 2 discusses the original 90nm cache design done by Xiaoyin Yao as part of his PhD dissertation. The cache design presented in this thesis relies heavily on this cache architecture. Chapter 3 deals with the detailed design of the 45nm RHBD cache circuit design. Chapter 4 provides the layout details of the cache data and tag circuits. Chapter 5 discusses the flow used for laying out the whole cache using structured array methodology. Chapter 6 concludes.

Chapter 2. 90 NM CACHE DESIGN

The original cache architecture and design was done by Xiaoyin Yao as part of his PhD program [47]. It was implemented on the IBM 90nm CMOS bulk process. He made use of specially designed logic cells to radiation harden the design. The following sections give an overview of the architecture and design details of the 90nm cache. This architecture would be used for the new design with some significant circuit modifications.

2.1 Introduction to Cache

Memory system follows a hierarchy as shown in Figure 2.1. As we go up the hierarchy, the memory is generally smaller and faster. These higher level memories normally reside on chip while the main and disk memories reside outside of the processor system.

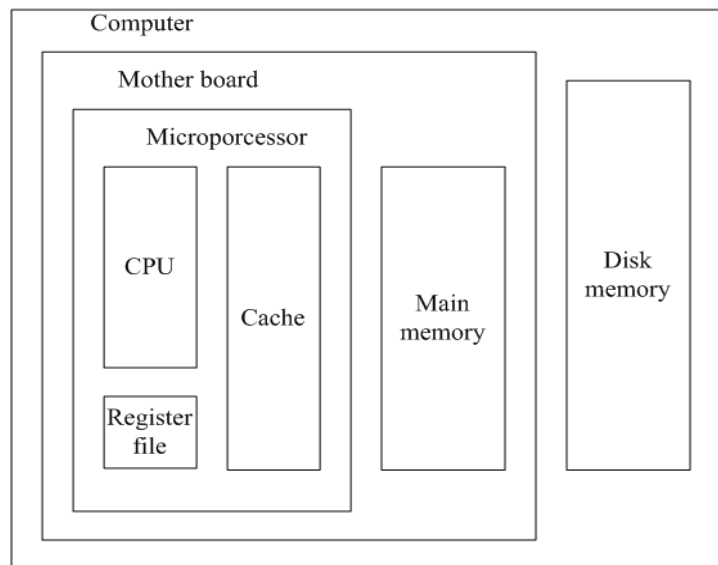


Figure 2.1 Memory Hierarchy

Cache works on the principle of locality. Consider the spatial locality and temporal locality of memory accesses. Spatial locality states that data whose addresses are near one another tend to be referenced close together in time. Temporal locality means that data recently accessed would be more likely to be accessed in the near future [48]. In more simple terms, the principle of locality means that same groups of data or instructions that were accessed recently are most likely to be reused soon.

Cache is used to store most recently used data and instructions. If the cache contains the data or instruction that the CPU is going to execute, the access is a hit. Otherwise, it results in a miss even though not every cycle accesses the cache. Hit rate defines the percentage of the total CPU cycles in which an access results in a hit. The miss rate is “the remaining percentage of CPU cycles” [49]. The miss rate (or hit rate) is an important measure of cache design. It is affected by the cache size and the write policies. There are two types of write policies: Write-through and write-back. In the write-through policy the information is written to the cache as well as to the lower memory. The write-back policy requires the information is written to cache only. Therefore, there is only one copy of the modified data. In the former case, there are two copies of data, one in the cache and the other in the lower-level memory. This is important when designing a rad-hard cache. For the cache with write-through policy, it only needs to detect a radiation induced error [47]. The correct data can be copied from the

lower-level memory. However, the write-back cache needs to correct the data by itself making it less ideal to function as a radiation hardened cache.

Cache can be used as either data cache or instruction cache. Data cache is used to fetch and store data, while instruction cache just fetches instructions. Data cache can be loaded from a lower level memory and store values which may be changed by the CPU. Instruction cache can only load from lower level memory, since instructions are not supposed to be changed.

If a miss happens during a read, a block in the cache needs to be replaced with the desired block from the lower level memory. There are many strategies employed to decide which cache block will be replaced. The block can be randomly selected or the least recently used. In this work, the replacement algorithm is based on the least recently filled (LRF) line. This algorithm is consistent with the standard MIPS architecture, which was chosen for the design. In addition, it is simpler to implement LRF [47].

Cache organization can be classified into 3 main categories: direct mapped, fully associative and set associative. They decide the block replacement policy in a cache. In a directly mapped cache, each block can only be mapped to one place. A fully associative cache can map a block anywhere. A set associative cache has restricted places or sets where the blocks can be placed. For example if there are n blocks in a set, the cache is n -way set associative.

If a miss happens during a write, there are two common policies that are employed: write allocate or no-write allocate. Write allocate caches replace the

missed block and perform a write again. No-write allocates caches only modify the block in the lower level memory.

The cache performance is a major factor that affects the overall throughput of a microprocessor. It is measured in terms of average memory access time:

$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

where hit time is the time taken for a hit in a cache, miss rate is the percentage that a miss happens, miss penalty is the additional time needed when a miss happens. Average memory access time can be measured either in absolute time or in number of clock cycles.

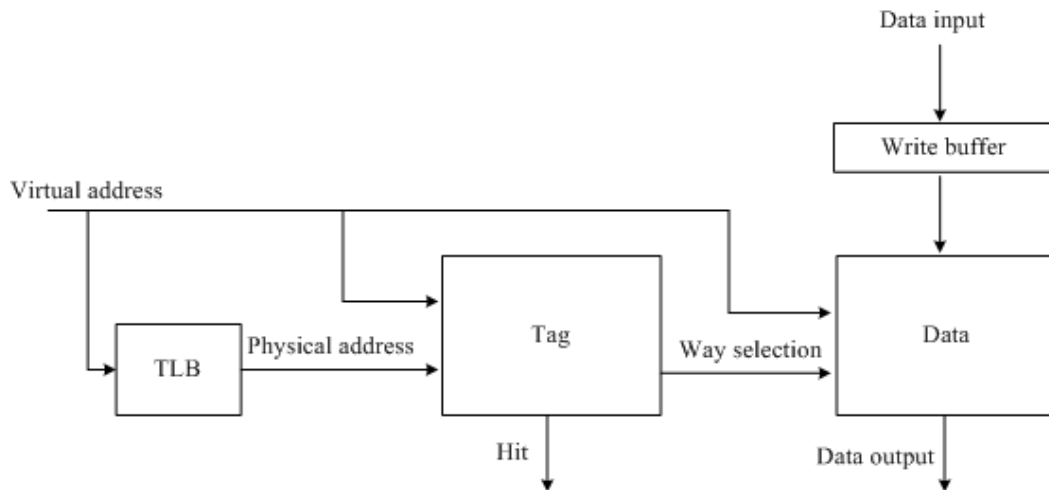


Figure 2.2 A high level example of the Cache structure [47].

Cache performance depends greatly on and can be improved by reducing the average memory access time. Some methods to reduce the access time are by reducing the hit time, miss rate or miss penalty. The cache is split into tag and data arrays so they can be addressed independently. This speeds up the hit

generation and data read out. A write to the cache is usually slower than the read because the tag must be read before writing the data. A write buffer which acts like a prefetch is often used, so that the write operation is pipelined. A translation look-aside buffer (TLB) is added to convert virtual address to physical address. Figure 2.2 shows a cache diagram with TLB, tag, data and write buffer.

Reducing the cache miss rate has been one of the greater focuses of cache research. The primary methods of reducing miss rates are by using larger blocks, higher associativity, pseudo-associativity, hardware prefetching, prefetching instructions, or compiler optimizations [48]. Figure 2.3 shows the miss rate for different cache size and associativity.

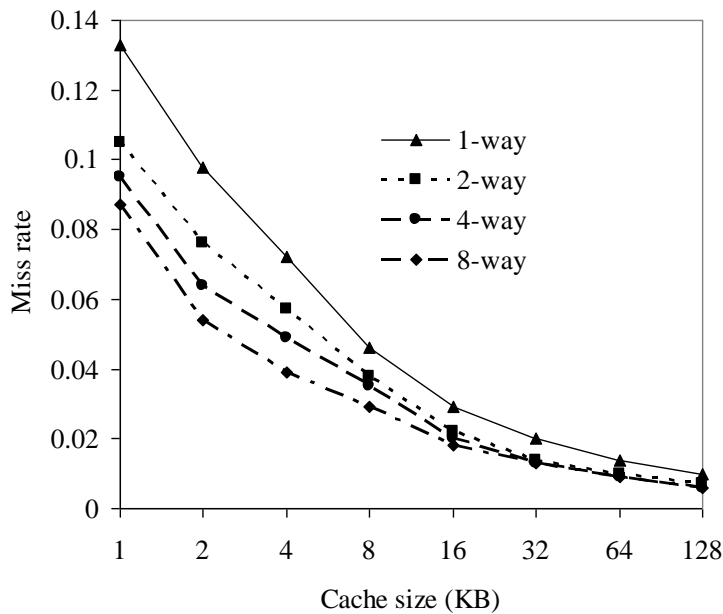


Figure 2.3 Miss rate for vs cache size and set associativity. (After [50])

The miss penalty can be reduced by giving priority to read misses over writes, using sub-block placement, restarting early and sending critical data first, implementing nonblocking caches to reduce stalls on misses, and adding a second-level cache [48].

Table 2.1 MAJOR INTEL[®] MICROPROCESSORS AND THE CACHE SIZE

Processor	Clock Speed(s)	Introduced Date	Process	Cache
Intel [®] Atom [®] 230	1.60 GHz	Jun-08	45 nm	512 KB L2
Intel [®] Xeon [®]	2.80-3.60 GHz	Feb-05	90 nm	2048 KB L2
Intel [®] Pentium [®] 4	2-2.60 GHz	Aug-01	0.13 μ m	512 KB L2
Intel [®] Pentium [®] III	500 M-1.13 GHz	Oct-99	0.18 μ m	256 KB L2
Intel [®] Pentium [®]	200-75 MHz	Mar-94	0.6 /0.35 μ m	8 KB

The above techniques no doubt improve the cache performance, but they add quite a lot of hardware and software complexity to the cache design. This

leads to more silicon area and power. Therefore, tradeoffs have to be made while planning to implement the above techniques.

Table 2.1 (after [51]) lists major Intel[®] microprocessors. The clock speed, introduced date, process node and cache sizes are included.

2.2 RHBD Cache Design Requirements, Assumptions and Approaches

This cache is intended to be used in a high performance radiation hardened microprocessor. The target operating frequency for the 90nm cache was 1GHz and the cache maximum operating power is less than 300 mW. The design should be able to work as either instruction or data cache. Thus a data cache is designed and some features are disabled when used as an instruction cache. There are two fundamental differences between the data and instruction caches. First, the instruction cache can only be loaded from a lower level memory. The data cache can be loaded or written by instructions executing in the CPU. Second, the minimum operation unit in the instruction cache is a word, while in the data cache it is a byte. Thus the data cache design is a superset and the work focuses on the same.

Two assumptions can be made about the radiation environment based on the underlying statistics [47]. First, when an error is caused by a SEE, a localized region (less than $25 \times 25 \mu\text{m}^2$ across) is affected for a very short time (less than 10 ns). Second, the probability that another error will immediately follow a SEE-induced error is small. The time duration between consecutive strikes could be many thousand nanoseconds. This applies to the space radiation environment near

the Earth. In GEO orbit, the average interval time between particles that hit the cache and has a LET of greater than 1 MeV-cm²/mg is 3×10^{12} ns [47]. Two major RHBD techniques were used in the 90nm cache design: Annular NMOS transistors and guard rings around NMOS transistors to achieve TID hardness. A detect-invalidation-reload scheme is used to achieve SEE hardness.

2.3 Cache Design and Architecture

The capacity of the cache is 16 KB and it is 4 way set associative, using write-through and no-write allocate policies. This size and set associativity gives a reasonable hit rate. The write-through policy is necessary to implement the detect-invalidate-reload scheme. Every eight data bits have one associated parity bit. The cache has 1024 cache lines, with 16 bytes of data in each cache line. The cache is virtually tagged and physically indexed. A cache line is the smallest division of a cache memory for which there is a distinct tag address. In this cache, a line consists of four words (with four bytes in each word) that have the same set address and way. Thus there are 256 cache lines (or 4 KB of data) per way.

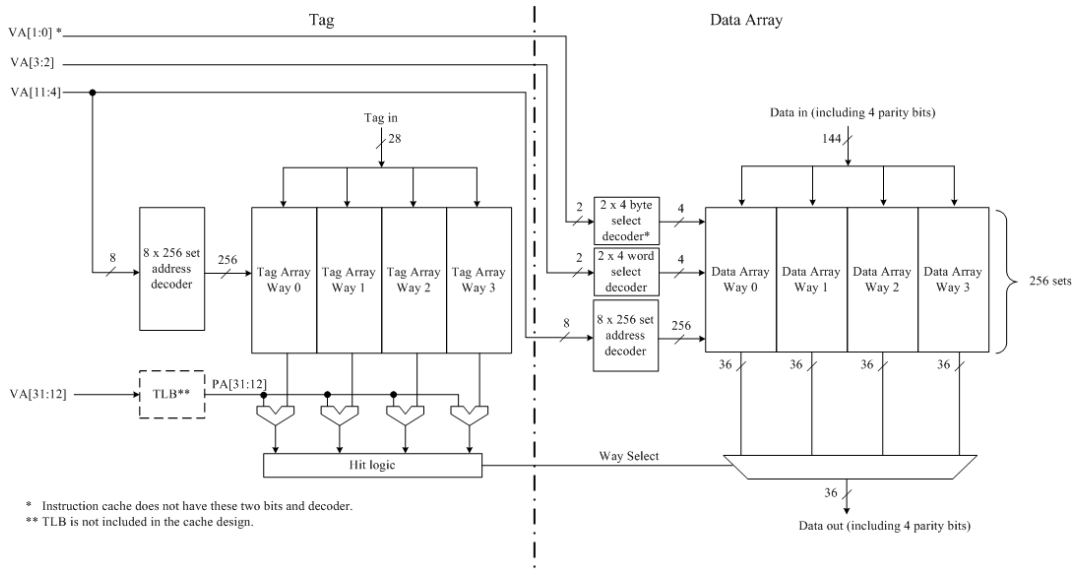


Figure 2.4 The basic diagram of cache. (After[47])

The most significant 20 bits from the virtual address are fed into TLB which generates the physical address, which in turn is used for the tag comparison. Unmapped virtual address bits are sent to the tag and the data array decoders. The way hit signal from the tag selects one of the 4 ways of the data array to generate the final cache data read out.

The cache supports four operations: lookup, read, write, and global invalidation [47]. In a lookup operation, all 4 ways are read and the tag address is compared with the physical address from the translation lookaside buffer (TLB). If there is a match, the selected way is chosen and a hit signal is generated. Otherwise, a miss signal is generated. This supports an instruction fetch (instruction cache) or a load instruction (data cache). The read operation can read a specified set and way from the data array or tag. The minimum unit that can be

read from the data array is a word. In the tag, the tag address and parity bits from a certain set and way can be read out. Other tag status bits belonging to all other ways are read out at the same time. A write operation writes a specified set and way of the data array or tag. The minimum unit written is a byte. In a line fill, a whole line is written. Global cache invalidation invalidates the whole cache. This operation clears the dual redundant valid bits in every tag entry. This occurs when the cache is reset after a power up, reset or a SEE-induced error is detected.

2.3.1 Data Array and Tag Array

This cache design achieves high performance and low power by using simple yet efficient design approaches. The design does not make use of the traditional sense amplifiers on the read paths. There are no column multiplexers due to the BL circuit architecture and word line (WL) arrangement. To achieve high performance, both data array and tag array use dynamic circuits to precharge and read their SRAM arrays. Gated clocks are used in the banks to reduce active power.

The floor plan of the cache is shown in Figure 2.5. The data array is divided into two halves, with the most significant 18 bits of each byte and parity protected word on the left and the least significant 18 bits on the right with the tag in the middle. This reduces the wire length of the way select signals and thus helps reduce power and speed up the cache operation.

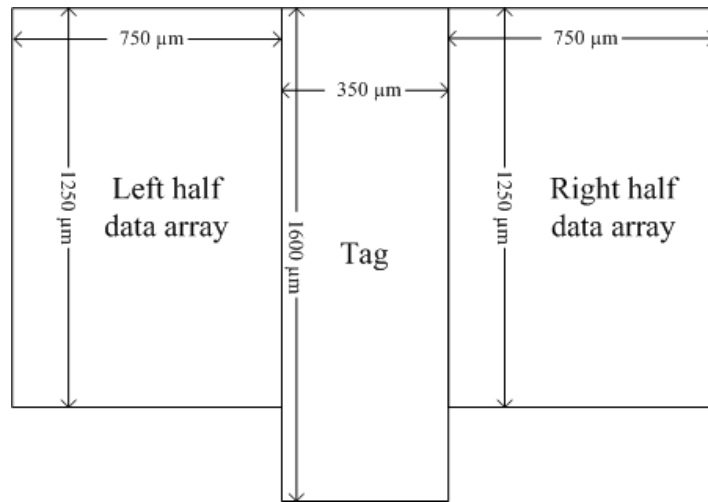


Figure 2.5 Floor plan of cache. (After [47])

The left or right half data array consists of 4 words. Each word comprises of 4 banks, consisting of a top and bottom sub-bank. A sub-bank comprises 32 rows and 72 columns of SRAM cells. In the data array, there are a total 16 KB of data. Since the smallest writable unit is a byte, it must be parity protected. Therefore, there are a total of 147,456 total bits of storage in the data array.

In each cell row, there are two word lines (WL) and two ways. In each WL of each way, there are two bytes. Each byte is composed of eight data bits and one parity bit. Figure 2.6 shows major circuits in the data array [47]. The data array fills in a single cycle, since all the words can be written simultaneously by activating the four banks at once.

There are two reasons why two WLs and two ways were implemented in a row. The first is to avoid column multiplexers which are usually used in a cache design. Therefore the BL development path is simplified, which helps the cache achieve high speeds. The second reason is that by interleaving bits belonging to

two bytes, two WLs and two ways, the bits belonging to the same parity group (same byte, WL and way) are separated by a distance equal to the width of seven SRAM cells, rather than three. This helps obtain SEE hardness against multi-bit upsets since a charge track would have to span 8 cells to upset two bits protected by the same parity bit.

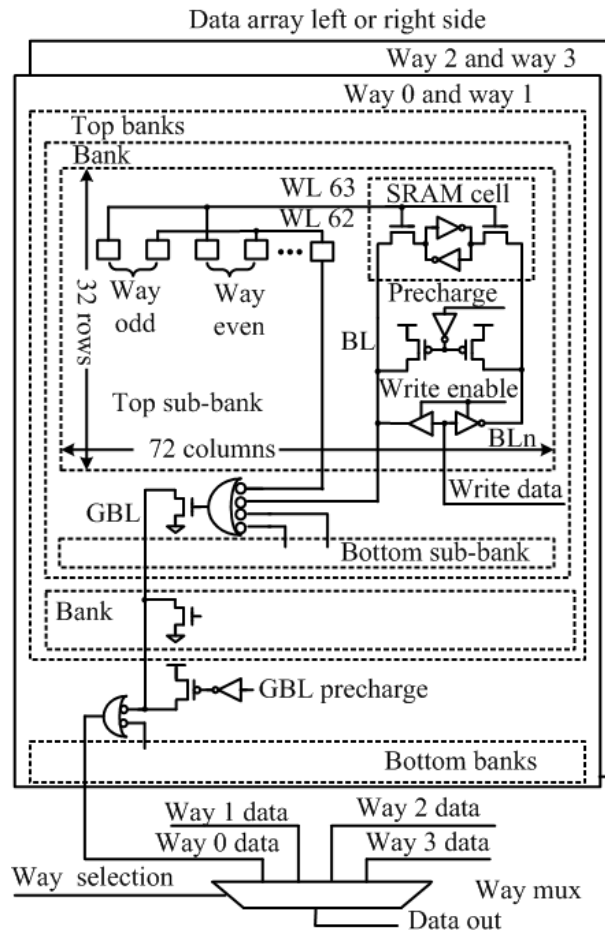


Figure 2.6 Basic diagram of data array (after [47]).

The precharge circuits for the bit lines (BL), write drivers and sense amplifiers are placed in the center of each of the half data array. The central location provides the shortest wires, again to reduce power and wire RC delay.

Redundant match lines are generated, and one of each is routed to the left or right halves respectively.

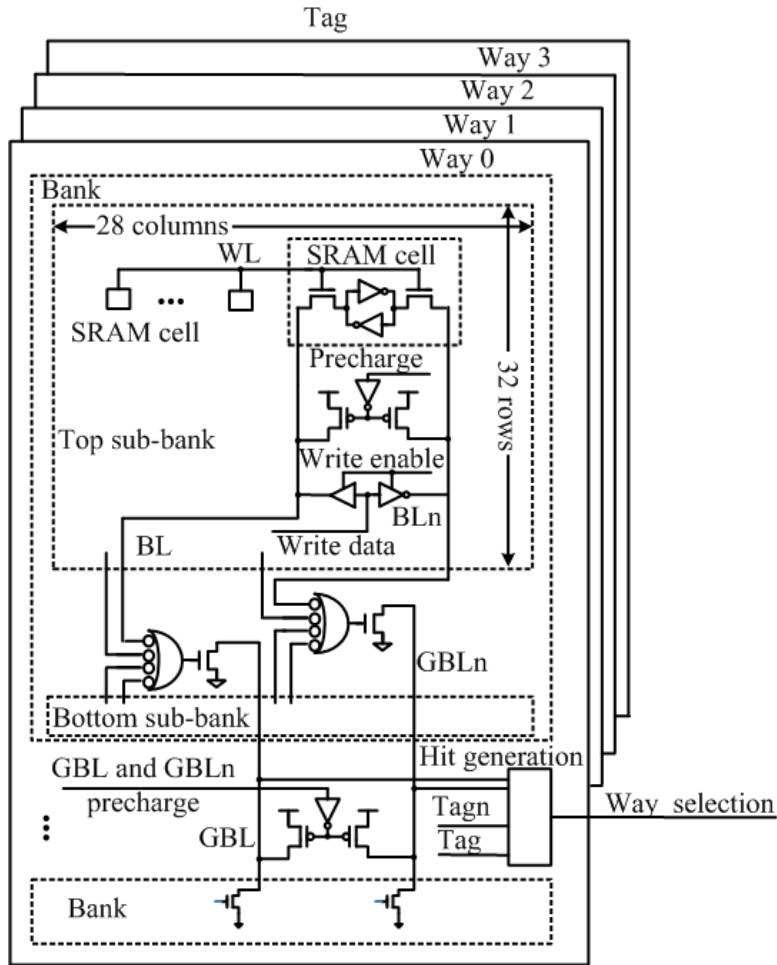


Figure 2.7 Basic diagram of tag (after [47]).

The tag array is composed of four ways, with four banks in each way and two sub-banks (top and bottom sub-bank) in each bank. Each way has 2 sets of redundant tag comparison and hit generation logic, so that SETs on them can be detected. There are 32 rows and 28 columns in a sub-bank. The components of

each way are shown in Figure 2.7. The tag array uses interleaved layout. There are 24 tag bits and 4 parity bits. The tag bits consist of a lock bit, a least recently filled (LRF) bit, two valid bits, and 20 tag address bits. The lock bit indicates a line that should not be replaced, making the cache suitable for real-time systems. The LRF bit indicates the least recently filled line. It is used to determine which line to replace after a cache miss. The valid bits are dual redundant and indicate the corresponding cache line is valid. For an invalid cache line, if either valid bit is set by SEE, the other valid bit would still indicate the cache line to be invalid.

2.4 “Staticizing” the 90nm Cache Design

The original 90nm cache design was almost fully designed using dynamic logic. Peripheral circuits like the way hit logic and comparators used dynamic logic. The objective here was maximum speed for a radiation hardened design but there are downfalls to this type of a design. The first among these issues is the timing characterization and second a more difficult problem, which is portability. Dynamic logic while very fast can consume a lot of power and can also lead to potential functional failures as nodes are minimally driven. Careful design practices are required, like the inclusion of keepers, which are a must though there is a small performance/area/power loss. The next few sections describe the original dynamic design, the static design, and some important design practices. The logic that was staticized was the critical tag way hit logic and readout paths.

2.4.1 The Dynamic way hit design

The tag address which has been read from the SRAM cells are compared to the TLB input bits to generate the way select logic. The design is shown in Figure 2.8 and Figure 2.9. The tag data coming from the SRAM cells through the single ended 4 input NAND sense amplifiers from each bank in a tag way is wire ORed together by dynamic bus lines before they are compared in the comparator. The 20-bit mismatch signals generated are passed into a D1 type domino OR structure before converging into a NAND set dominant latch. The total delay from the clock rising edge to the cache final data read output is 930ps (after [47]).

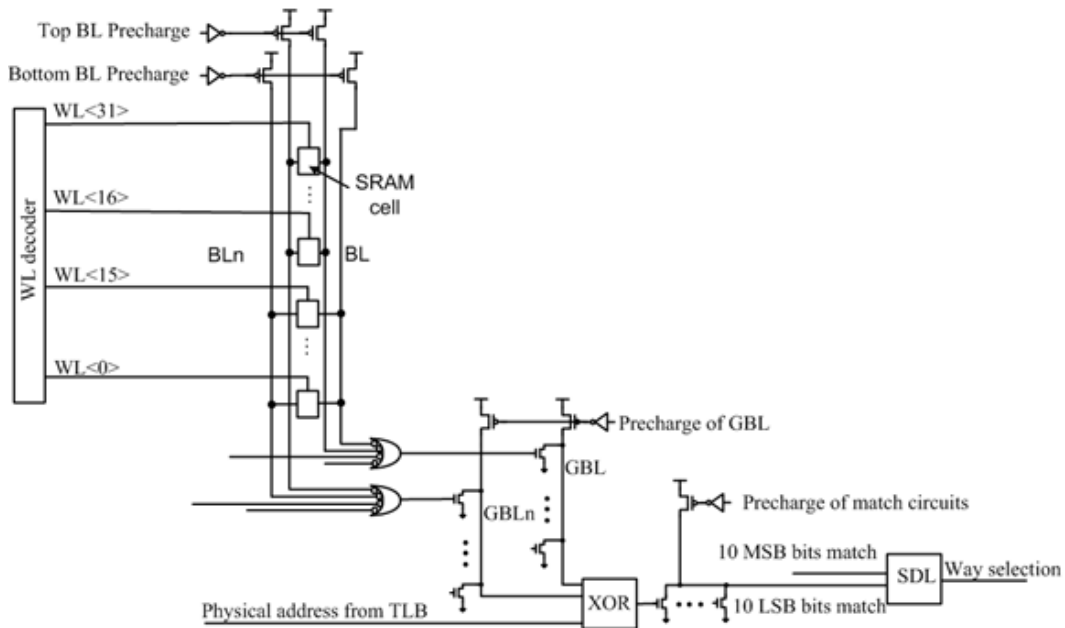


Figure 2.8 Overview of the Tag critical path. (After [47])

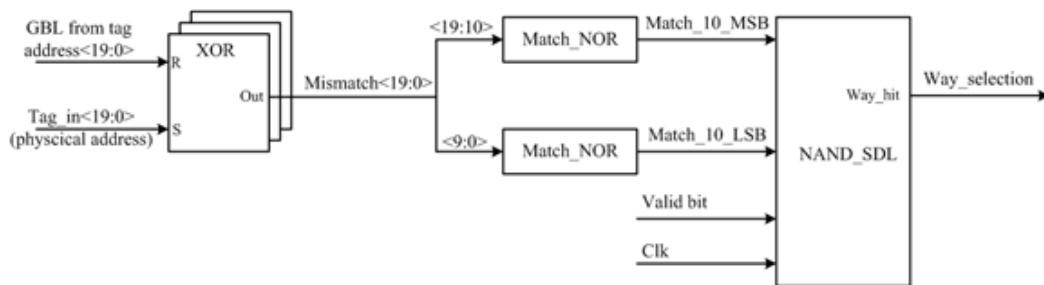


Figure 2.9 Hit Generation Circuits. (After [47])

2.4.2 The Static way hit design

The original dynamic logic based design was replaced with the more robust static CMOS based design. The schematic captured from the schematic editor of the entire modified way hit logic to the data read output is shown in Figure 2.10. A 4-input NOR gate combines the tag address readouts from the SRAM cells from the 4 banks of a tag way. The comparator used is a static CMOS 2 x 2 input OAI gate. CMOS 2 x 2 input OAI gate.

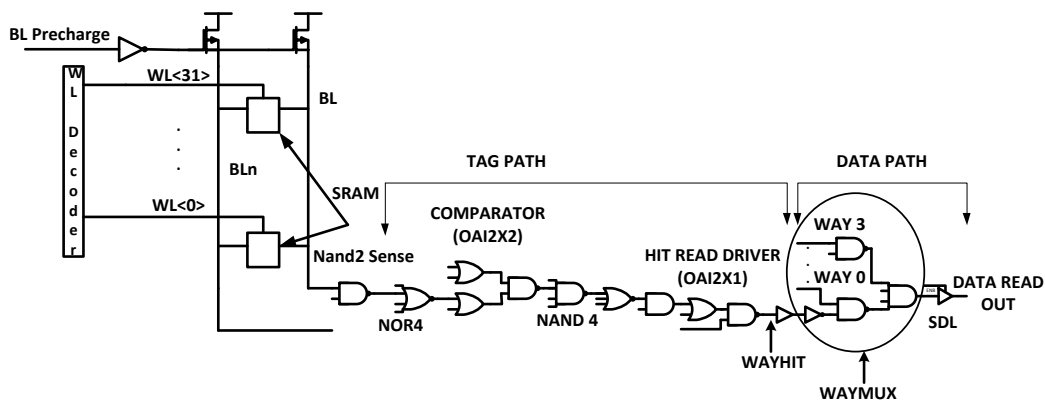


Figure 2.10 Static way hit schematic

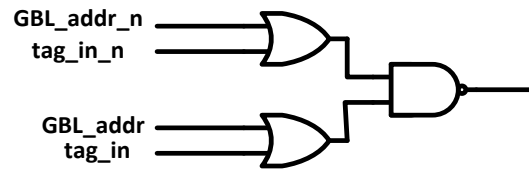


Figure 2.11 2x2 OAI Comparator static logic

The gate as shown in Figure 2.11 basically does an XNOR operation on the global address to generate the match signal. This signal propagates through a series of combinational gates to generate the way hit signal to the data banks. Figure 2.12 shows the simulation and the delays for the critical path. The delay when compared to the dynamic case is slightly greater but acceptable. This proves that a highly optimized static design can nearly equal its dynamic counterpart in terms of performance while the main advantage is that substantial dynamic power dissipation is saved.

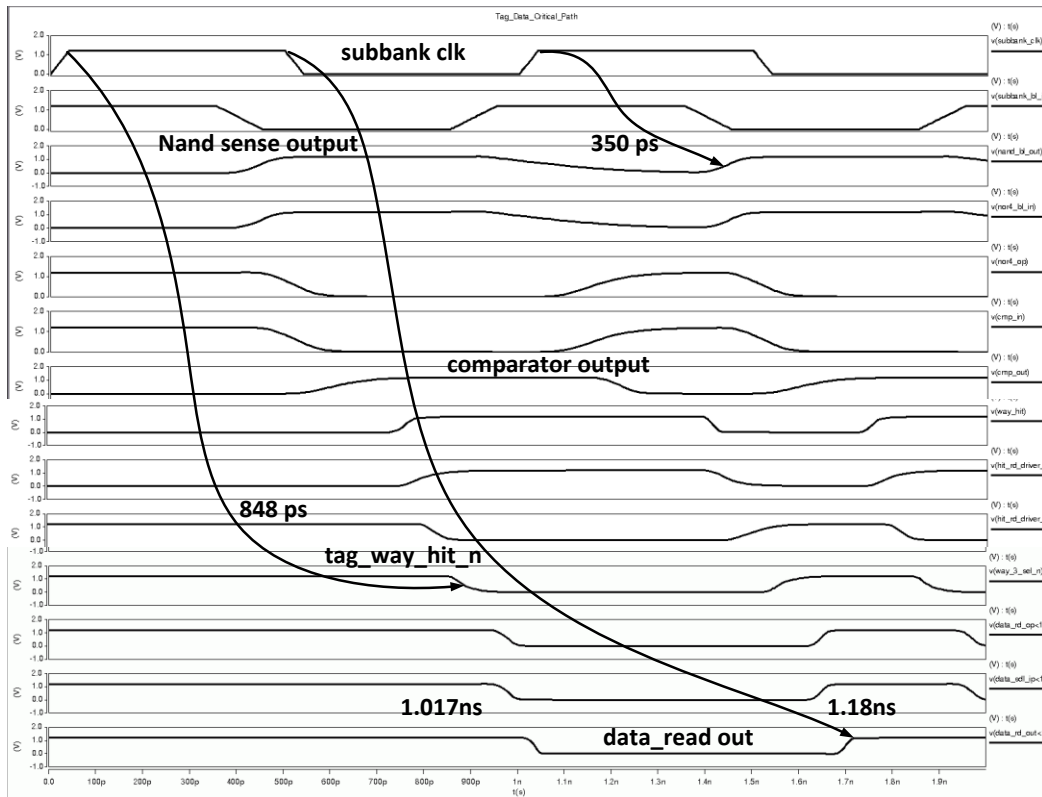


Figure 2.12 Simulation waveform for the cache critical path

The final data is read out from the data array through an inverter set dominant latch (SDL). The SDL provides the dynamic to static conversion. The bank discharging is the beginning of a dynamic path. Any static path that comes after a dynamic path will just let the dynamic circuit precharge edge propagate through. The SDL is used in such a way that it only allows the precharge to ripple through only when during the evaluate phase, otherwise it retains the previous value. Care is taken such that during precharge the SDL has a logic high at its input, otherwise the precharge can corrupt the data. The schematic as well as the simulation waveform is shown in Figure 2.13 and Figure 2.14 respectively.

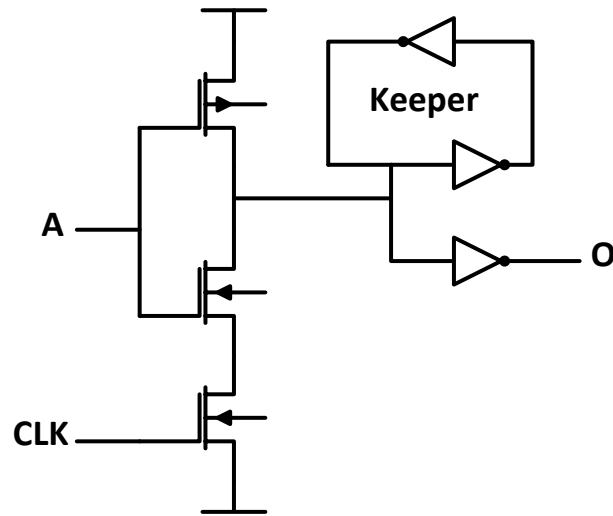


Figure 2.13 Schematic of an Inverter SDL

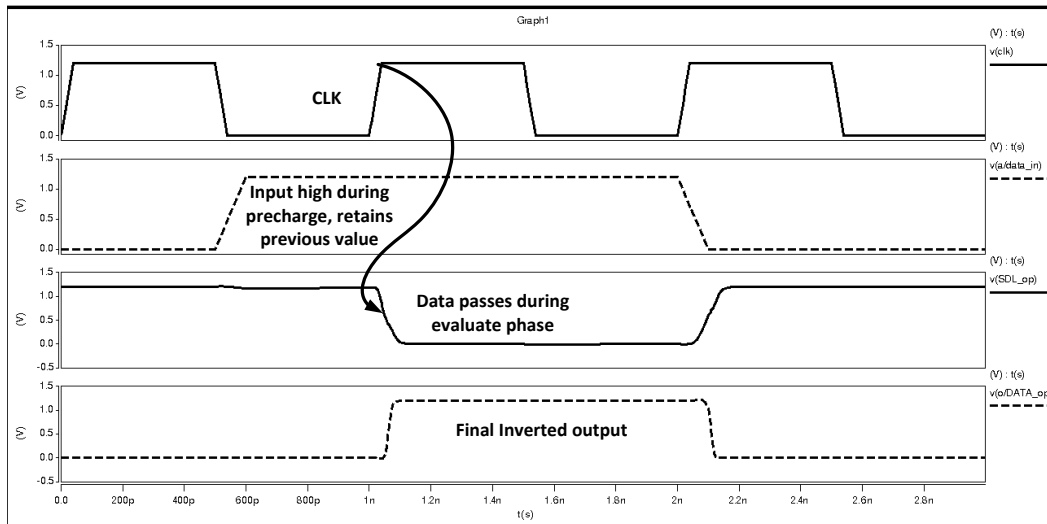


Figure 2.14 Simulation snapshot for an inverter SDL

The output inverter is a must while designing an SDL owing to the fact that a bare exposed node of a SDL can be susceptible to noise issues which can lead to functional failures. The keeper gives a good robust dynamic logic design

but the keeper sizing should be small else the pull down would have to fight with the keeper during the evaluate phase.

2.5 45nm RHBD Cache Design

The 90nm RHBD Cache has been proven to be very hard in SEE beam testing and TID testing [34] [60]. It features a TID hardened SRAM cell as shown in Figure 2.15. The layout is shown in Figure 2.16. It uses annular devices.

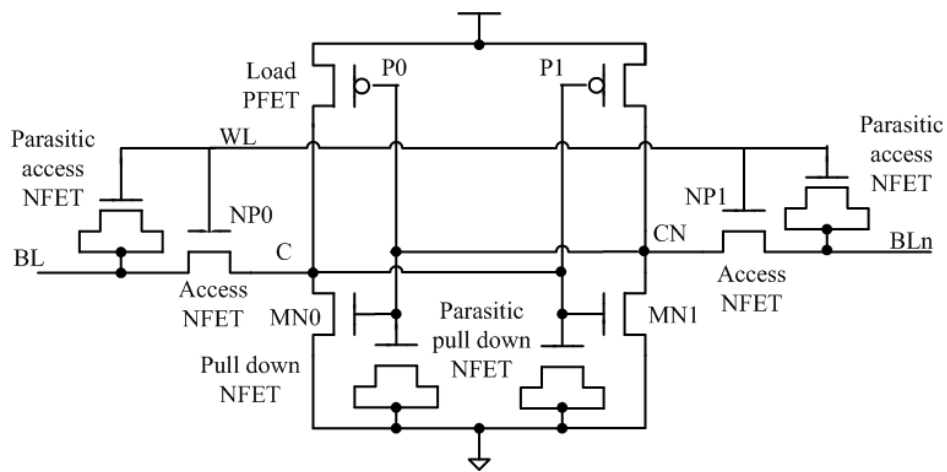


Figure 2.15 Modeling the parasitic neck loading due to the annular devices (after [47]).

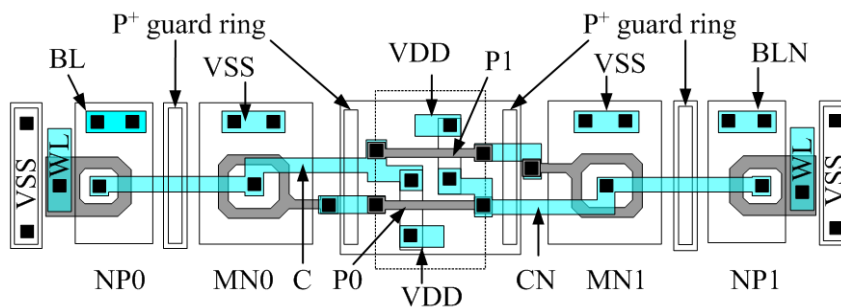


Figure 2.16 Layout of the NMOS-access RHBD SRAM cell (after [47]).

The TID hardening produces a large area and also requires the use of larger logic peripheral circuits. This slows down the cache operation. Additionally, as

mentioned, the entire design made extensive use of dynamic circuits which is very difficult to characterize and port to a different technology. The proposed cache design tries to address these issues by using a regular foundry cell from the IBM 45nm SOI process and by the use of static standard cell logic.

2.6 IBM 45nm SOI issues

Silicon on Insulator (SOI) has been an interesting area of research for decades. SOI is attractive since it has the potential for higher performance and lower power. The main difference between SOI and the conventional bulk process is that source, drain and body are surrounded by an insulating oxide rather than a conductive substrate. This eliminates most of the parasitic capacitance of the diffusion regions. Figure 2.17 and Figure 2.18 show a SOI inverter cross-section and an electron micrograph, respectively.

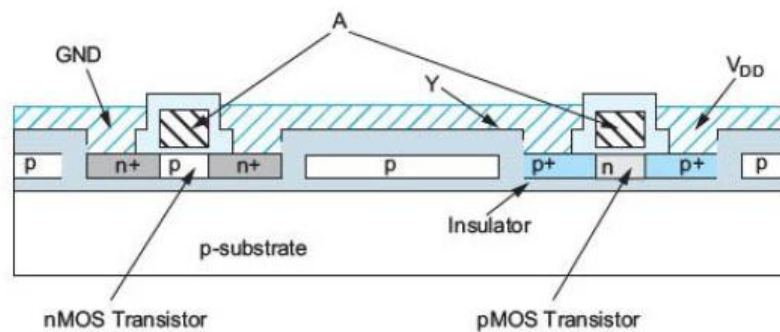


Figure 2.17 SOI Inverter cross-section (after [52]).

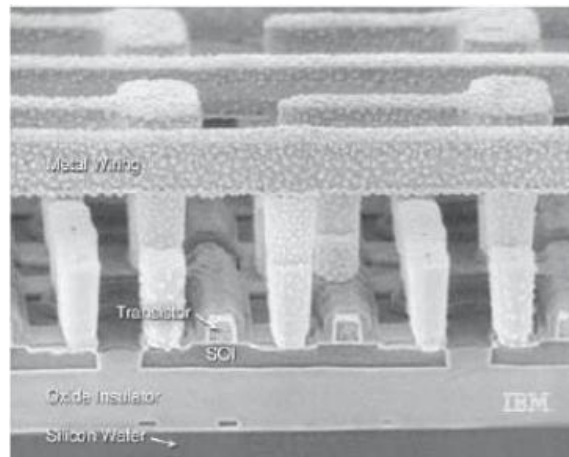


Figure 2.18 IBM SOI process electron microscope snap (courtesy of IBM).

However SOI comes with higher manufacturing cost and unusual device behavior which complicates circuit design. Using an insulator means that the body is floating as it is no longer tied to V_{DD} or V_{SS} through the well. Any change in the body will lead to a modulation of the threshold voltage (V_t). It can also lead to higher device temperature. SOI suffers from history effects whereby changes in the body voltage modulate the V_t , causing variable gate delays. Gate delay becomes a function of the switching history as the body voltage depends on whether the device was switching or staying idle. Figure 2.19 shows the charge paths to and from the floating body.

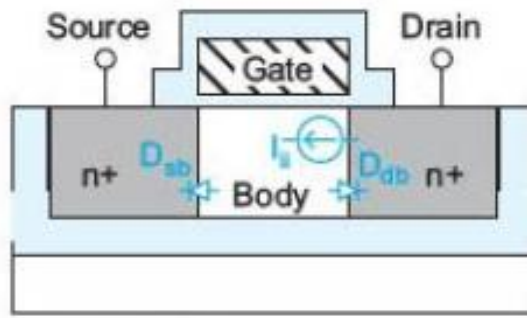


Figure 2.19 Charge paths in an SOI device (after [52]).

A reduced threshold voltage makes the gates faster but the uncertainty poses a challenge for circuit designers. These effects can be modeled and given sufficient guard band in a design, the critical timing variations can be taken care of. One can bound delays by subjecting the design to various conditions like initializing the body to V_{DD} for fastest performance or to V_{SS} for the slowest. History effect causes significant mismatches also between identical transistors. Another problem is due to pass transistor leakage. There is a parasitic bipolar transistor in each transistor as shown in Figure 2.20. In a bulk process the bulk is connected to supply whereas in SOI the body floats. If the source and drain are both held high for a long time while the gate is off, the base floats high due to diode leakage. When the source is pulled down low the parasitic transistor turns on. This results in a flow of current from drain to source even though the gate is off, causing leakage. This does not affect static circuits much, since one of the complementary gates is on. This can lead to incorrect functionality for pass gates as well as

dynamic logic. Therefore strong enough and carefully designed keepers are required for holding dynamic nodes steady.

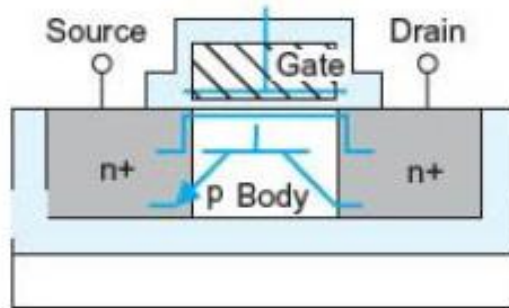


Figure 2.20 Parasitic bipolar transistor in an SOI device (after [52]).

Figure 2.21 shows the pass leakage example [52]. The dynamic node X is precharged high initially and the gate connected to this node is off. The source is high and discharges low producing a current stream which partially discharges the node X.

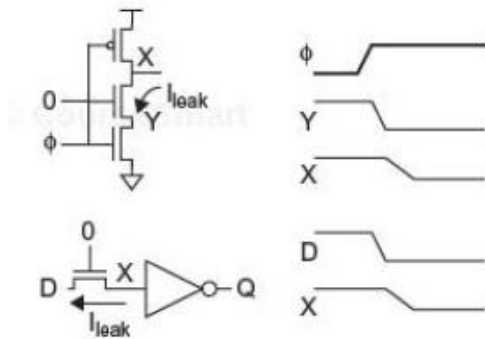


Figure 2.21 Pass gate leakage in dynamic gates (after [52]).

SOI can also lead to self-heating in devices. This is due to the insulating oxide which serves as a good thermal as well as an electrical insulator. The heat dissipated tends to stay inside the transistor rather than dissipate to the substrate. This may lead to slower operation and higher thermal dissipation in clock drivers or IO drivers. This leads to hotter chips. The issues described above are manageable, and if the designer takes care of all these issues during the design phase by including margins, the benefits of SOI can be reaped.

Chapter 3. **CIRCUIT DESIGN**

The goal of this work is to port a RHBD high performance cache design from a 90nm CMOS bulk process to a more recent 45nm SOI process. The technical difficulties presented by this type of logic was overcome by the implementation of many of the logic like the tag way hit and read out by using static logic standard cells from ARM. This cache design makes use of standard foundry 6T SRAM cells from the IBM 45nm SOI library unlike the 90nm one which used specially designed 6T SRAM cells to achieve radiation hardness. The decision to use normal cells was to achieve faster design time and a much smaller design which would not have been possible with the use of larger cells. On the downside TID radiation hardness decreases but SOI and the use of complex checking circuits and dual redundancy alleviates SEE. Section 3.1 lists the performance and radiation hardness requirements. The cache is planned to be fabricated on an IBM 45 nm SOI foundry process.

3.1 Cache Design Requirements

This cache will be used in a high performance radiation hardened microprocessor, which determines the performance requirements. The target operating frequency is more than 2 GHz and the cache maximum operating power is less than 300 mW.

3.2 Data Array Design Details

The data array needs to perform basic function of a memory unit, including write and read at 2GHz clock frequency. A number of circuit and

micro-architecture techniques are used to achieve radiation hardness. Subsequent sections describe the basic design features of the data array. Low power is another goal of this design.

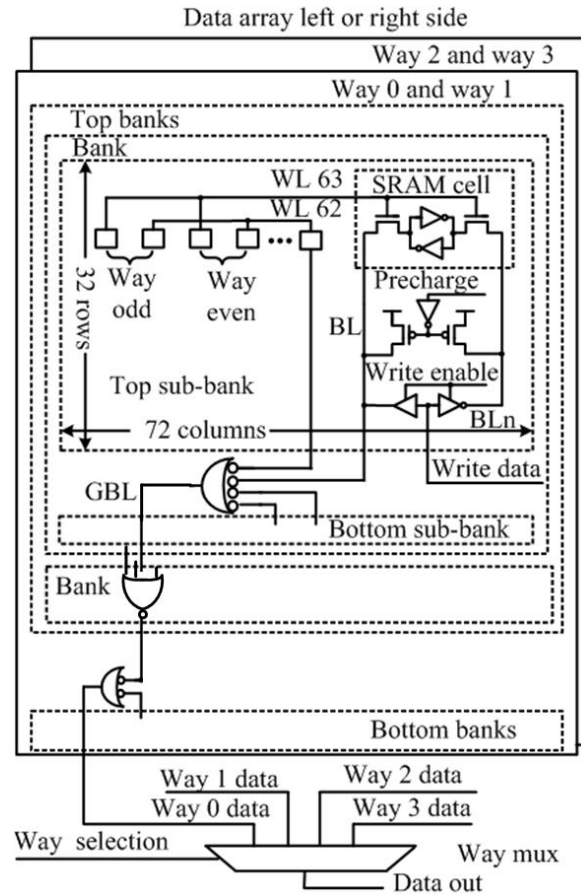


Figure 3.1 Basic diagram of data array.

The major component circuits in the data array include the WL decoder, SRAM cell, precharge circuitry for BL, write control, and the way multiplexer, as shown in Figure 3.1. During the precharge stage, the BL is pulled up to V_{DD} by PMOS transistor. When the clock is asserted and there is an operation to a particular bank, the WL decoder asserts one WL high. One SRAM cell in each

column or an entire row selected by the WL is selected for a read or a write. On a write operation, the BLs are driven by the write drivers and the selected SRAM cells are written. During a read, the BLs are driven by the selected SRAM cell. The read out value propagates through the NAND4 sense gate. The BL of the top four banks and bottom four banks are combined through a high-fan-in logical AND function implemented by standard cell logic. For each way, the result is the read out. These are passed to the way multiplexer. The way selection generated by the tag hit signals decide which way to read out to the cache output. The final read out data leaves the cache through a set dominant latch (SDL).

The WL decoder drives 64 WLs in a sub-bank. There are two WLs in each row, with 36 SRAM cells on each WL and is quite heavily loaded. There is no column multiplexer. Instead, a nand4-driver pull-down is used. Because of the cache architecture, only one of the four inputs of a nand4 gate at any given time asserts low. The other three inputs remain in the precharge state. This is because only one of banks belonging to the same way is active during a read or lookup. This results in fast operations.

3.2.1 SRAM Cell

The SRAM cell is the usual 6-transistor structure as shown in Figure 3.2.

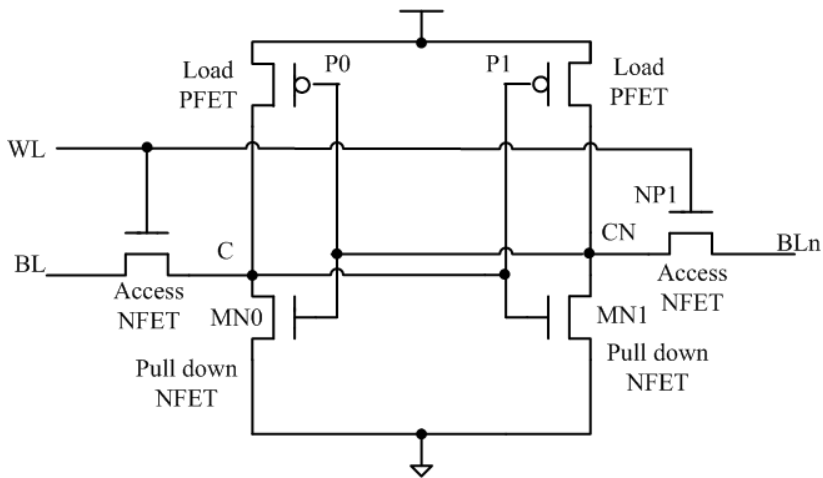


Figure 3.2 Schematic of the SRAM cell.

The fundamental SRAM cell used is an IBM foundry cell S462 which is used to build high performance memory arrays. The cell has an area of 0.462 μm^2 . The cell is written differentially, with one bit line (BL) and the other (BLn) at opposite potentials. It follows all SRAM special DRCs from the IBM 45nm SOI foundry process.

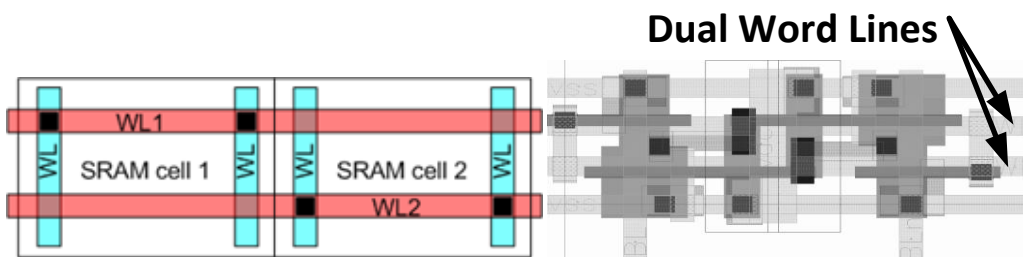


Figure 3.3 Two word lines are connected alternately to cells in one row.

Table 3.1 AREA COMPARISON BETWEEN THE SRAM CELLS FOR 45 NM PROCESS

PROCESS	IBM 45NM SOI
DEVICE	AREA (μM^2)
IBM S462	0.476
EQUIVALENT LOGIC CELL	0.753

lists the SRAM cell device sizes of all different SRAM cells in the 45 nm process used in the design. Note that the logic cell consumes twice the area of a foundry SRAM. The larger access and pull down devices in a SRAM cell provide high current to drive the high BL capacitance at high slew rates. The foundry cell is tall enough to allow two word lines to be interleaved in the layout—parity bits thus have twice the interleaving that would otherwise be not possible, i.e., have greater critical spacing. Also, this halves the access power since only every alternate cell is accessed when one WL is asserted (see Figure 3.3) [47].

3.2.2 Word Line (WL) Decoder

A static WL decoder is used in this cache design as opposed to the original dynamic decoder used in the previous design. This also helps to lower power dissipation. The decoder has 3 parts, namely a predecode stage, a postdecode stage and a clock and stage to produce the final wordline enable. The

combinations of the outputs from these predecoders generate 64 WLs, as shown in Figure 3.4.

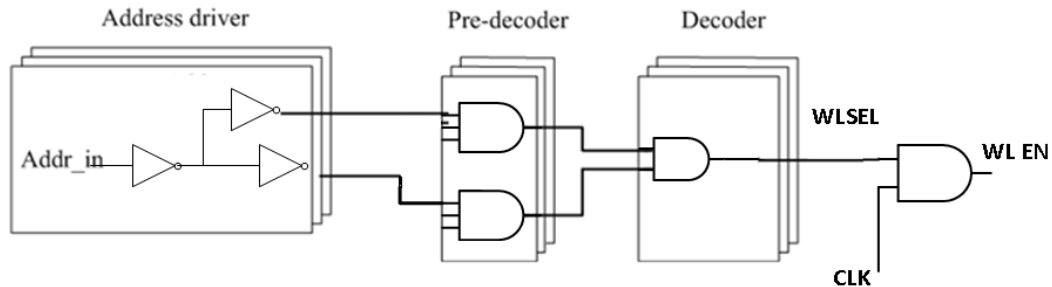


Figure 3.4 WL decoder.

The WL addresses are qualified by the sub-bank clock which is active only in selected banks. There is no active power in unselected sub-banks. The overall cache power dissipation is thus reduced by using this scheme. It toggles only when the sub-bank is active. If a sub-bank is not active, it is kept in the precharge stage. The decoder outputs (wordline selects) are available before the clock arrives. Decoding occurs in the negative clock phase (clock = low). A negative latch passes the addresses arriving during the negative phase of the clock. Decoding takes place in the negative phase of the clock and wordline selects are available for ANDing with the clock when the positive edge arrives.

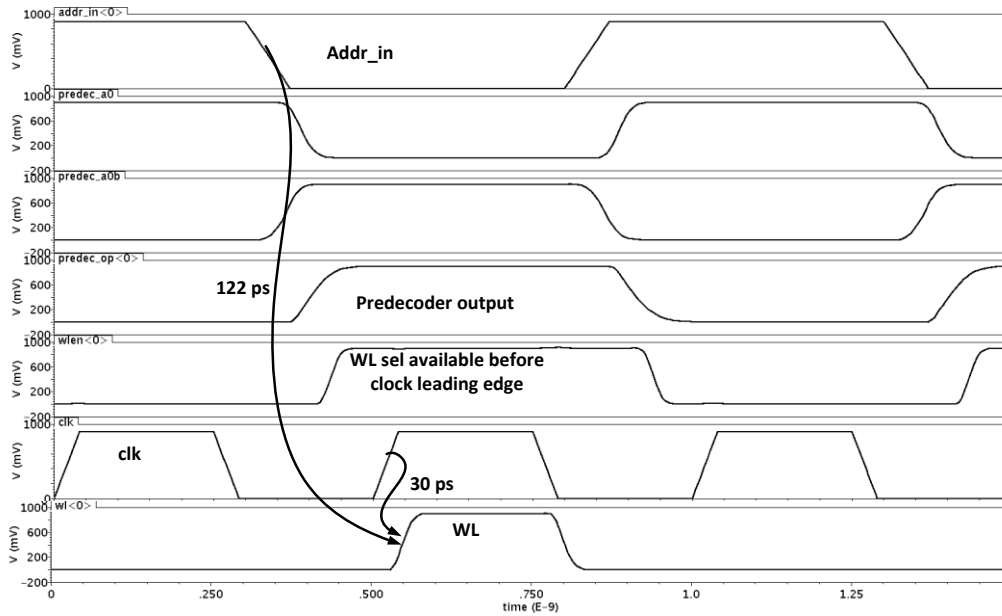


Figure 3.5 WL decoder simulation results.

Figure 3.5 is the simulation of the WL decoder. The longest delay from the rising edge of clk to WL is about 30 ps, at room temperature and at the typical process corner.

3.2.3 Write and Precharge Circuitry

As shown in Figure 3.6 [47], the BL and BLn are precharged by PMOS transistors. The write data is driven by static buffers and inverters. For writes, both BL and BLn are driven by tristate inverters controlled by the write enable signal (WREN). BL and BLn are precharged to V_{DD} by PMOS transistors. To ensure the WL, precharge and WREN work correctly together, a self-timed circuit is used for the write and precharge control, respectively, as illustrated in Figure 3.7 [47]. The timing is derived directly from the sub-bank clock.

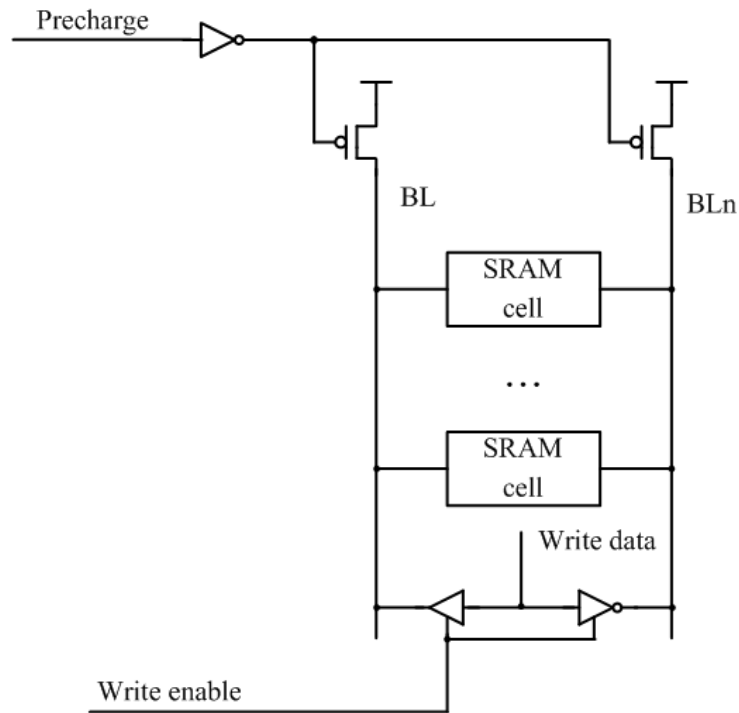


Figure 3.6 Write and precharge circuitry. (After [47])

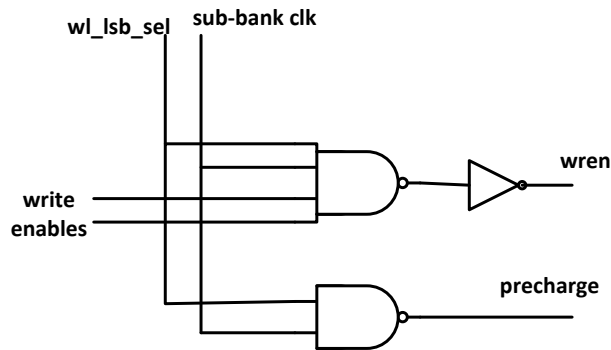


Figure 3.7 Self-timed circuits are used in WREN and precharge drivers.

Both WREN and precharge are generated from the sub-bank clock. There are two inversions in the WREN driver and one inversion in the precharge driver. Although the loads on the WREN, precharge and WL are different, the driver

sizes have been adjusted to make the delay on these paths similar to each other.

Figure 3.8 illustrates the simulation results on these paths.

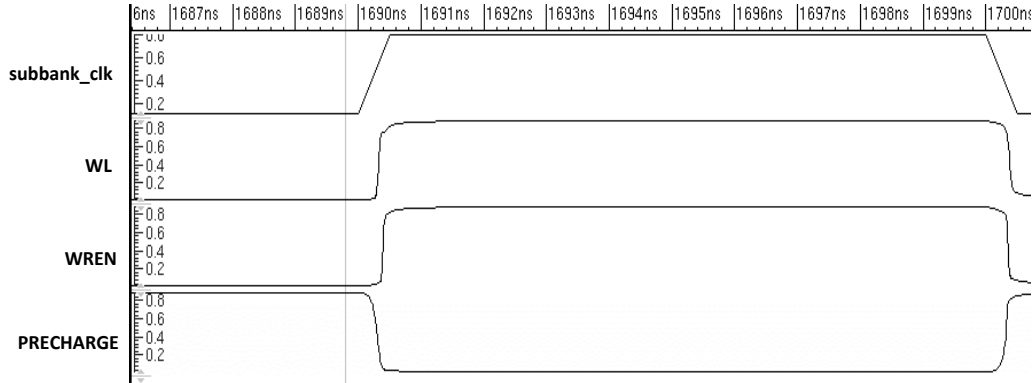


Figure 3.8 Simulation waveform of WL, WREN and precharge.

3.2.4 Techniques to Achieve SEE Radiation Hardness

This section focuses on the various techniques used in the data array to achieve radiation hardness by design. The detect-invalidation-reload scheme is used for cache SEE protection. It involves both architectural and circuit level techniques. The cache is write-through and we always have a copy of the data stored in the lower level memory. This scheme exploits this very nature of the cache. The whole cache is invalidated when there is a SEE-induced error. The processor is restored to a valid state over time after the SEE. The valid architectural state may be retrieved from an external protected memory source that resides outside of the processor system. This design makes use of parity protection, interleaved layouts, dual redundancy and special error checking circuits to mitigate SEE errors and MBU.

Layout spacing plays a major role in reducing errors due to MBU which cannot be detected using parity bits. SRAM bits from different parity groups are interleaved with each other. The bits belonging to the same group are separated by a width equivalent to 7 SRAM cells. This reduces the chances of MBU affecting all the bits in a single group. Within each group single-bit-upsets can be detected by their respective parity bits.

This cache design employs clever error detection circuits to monitor the cache operation in case of a SEE. They are implemented in the architecture such that in the event of a SEE, they report the errors before the data is read out and used in the pipeline. The error circuits can themselves be victims of SEE. In such a case, there might be false alarms. The entire cache is invalidated if such a scenario occurs. Figure 3.9 shows the detection circuits implemented in the data array. All inputs to the data array are dual redundant at the boundary. The dual redundant signals are compared before the signals enter the data array to check for any errors. A list of possible errors and corresponding detection techniques are described below.

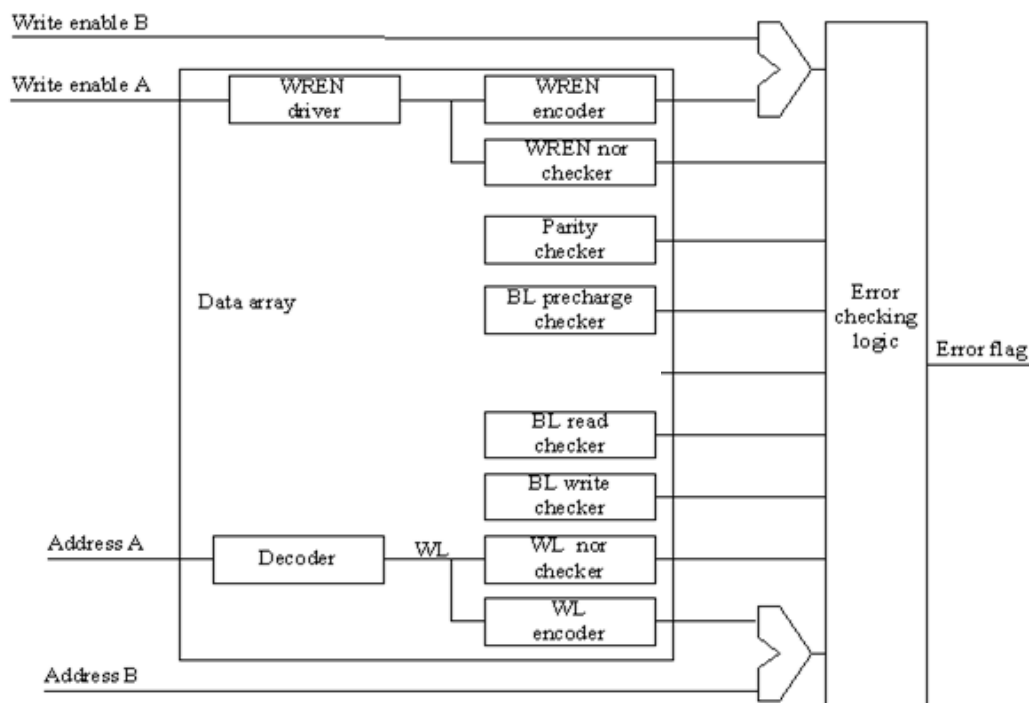


Figure 3.9 The error detection scheme of data array. (After [47])

When one or more SRAM bits are flipped, owing to the layout spacing between bits in the same group, the parity checks catch them as they all get converted to single-bit errors.

In another case of SEE error, the write enable signal may get asserted incorrectly. The write enable signals are made dual redundant, indicated as write enable A and B in Figure 3.9. The path A is used in the data array. The write enable driver controls the write circuits for individual SRAM columns. The WRENs are encoded back and compared to the WREN from path B. An error is flagged if there is a comparison mismatch. There can be a case of an undetected

error that occurs in the write circuit of a single SRAM column. The parity checker will duly catch this error, when the data is read out since it is a single-bit upset.

The address driver or WL decoder can give incorrect outputs in case of a SET. The address input protection also involves similar measures as the write enables. There are dual redundant paths for address inputs. One path is fed into the data array WL decoder. An encoder within each bank regenerates the address signals from WLs. If there is an error in the address decoding path, there will be a mismatch between the regenerated address and the redundant address flagging an error.

The BL precharge can be asserted or de-asserted incorrectly, due to an SET. There are BL precharge checkers to detect any error that occurs in the precharge control logic.

The BLs or BLns could be wrong during a write or read cycle due to a control signal error or WL SET induced error. A BL read-checker can report errors if either BL or BLn fail to develop fully during a read phase. The BL write checker reports an error based on the event that a special dummy SRAM cell is not written correctly during a write phase.

Finally all the errors signals from the checking circuits are ORed together in the error checking logic outside of the cache. Any error from the data array will lead to the cache being invalidated and is reported for test purposes. The errors can also be logged, so that we can determine their relative frequency and cross-

section vs. LET. The following sections describe the design details of each checker.

In the data array, the detection circuits comprise roughly about 15% of the total area. The detection circuit impact is similar in the tag array. Therefore, it can be concluded that the overall cache area penalty of using the detection circuits is around 15%. This is a small yet significant area penalty for achieving radiation hardness for the design.

3.2.4.1 Error Checking Circuitry

The checkers should be designed such that they catch any error that occurs instantly and not just the one that occurs during a clock edge. Static checking circuits thus become inappropriate for this purpose. All of the error checkers are based on the same dynamic logic design. The basic dynamic error checker is illustrated in Figure 3.10. It detects an error when the checked signal is high (checking window). At all other times, node A is precharged to V_{DD} . The error signal is combined with error signals from all the other banks, in the error control logic. Some error flags propagate into a latch or flip-flop directly because they may get updated with a new value in the next check cycle. The checking circuits implemented for the WLs, write enables, BL precharge, bit line read and write follow the same approach.

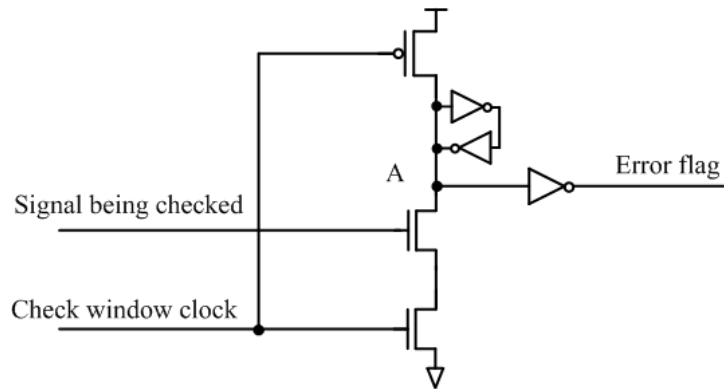


Figure 3.10 The basic dynamic error checking circuit (after [42]).

3.2.4.1.1 WL Encoders and NOR Checker

A WL error during a write, can lead to a missed write operation or data being written to the wrong row. As the parity is stored at the same time, therefore the parity checker cannot detect this error, leading to MBU. This is commonly referred to as a silent data corruption which is fatal. If such an error happens to WL during a read, a wrong row or multiple rows can be read. The dynamic WL encoder regenerates the address based on WL(s) actually asserted. The regenerated address is compared with the input address to the WL decoder. A comparison mismatch indicates a WL error. The scheme and structure are shown in Figure 3.11 and Figure 3.12.

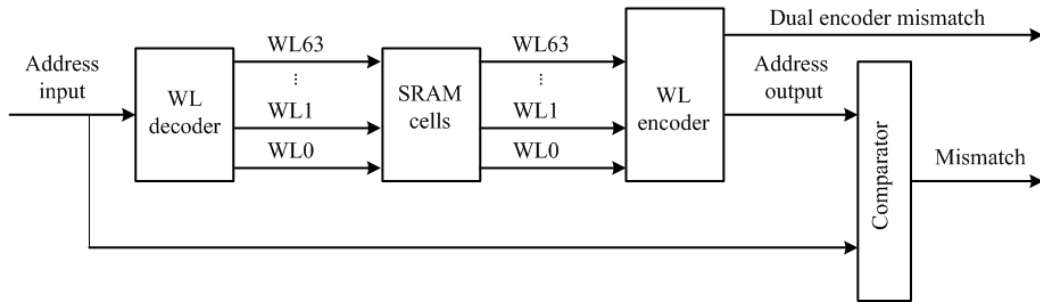


Figure 3.11 The WL encode scheme [47].

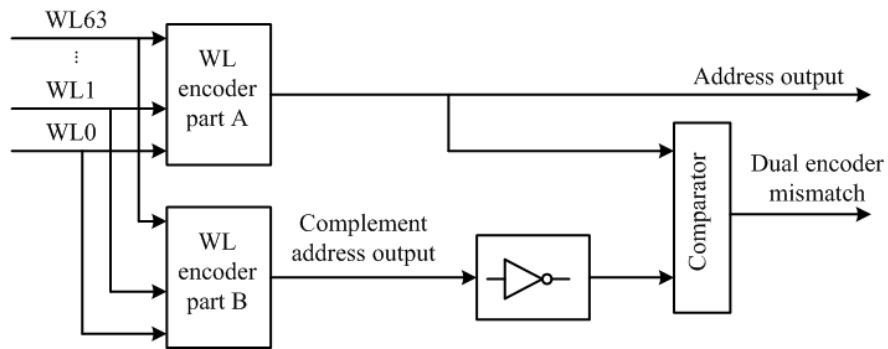
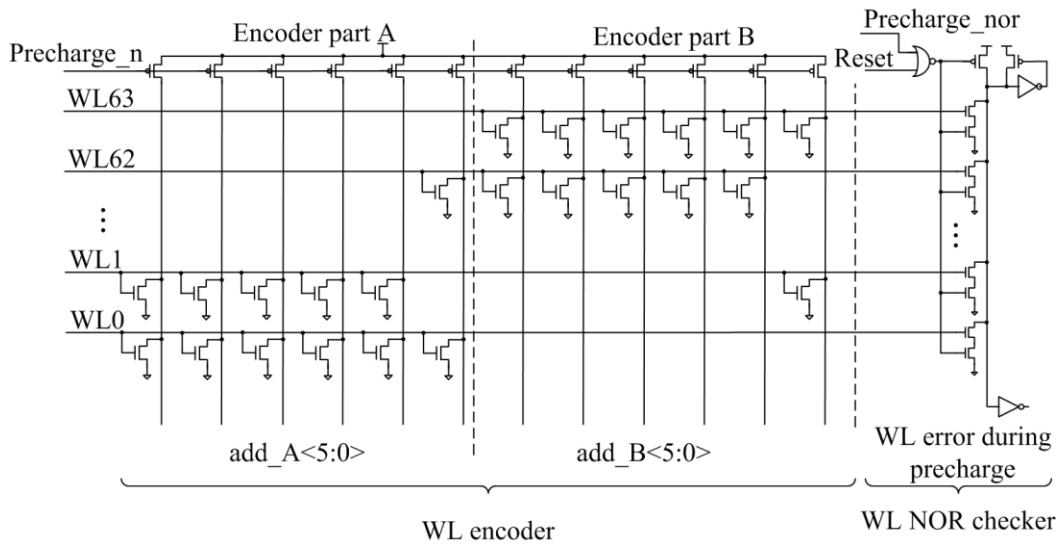


Figure 3.12 The WL encoder structure [47].

Incorrectly asserted WLs can also cause timing errors. A separate circuit, NOR checker, checks for WLs asserted during the precharge stage. The WL NOR checkers are active even in banks that are not active.



Note: there are keepers on each vertical line although they are not shown.

Figure 3.13 WL encoder and NOR checker (after [47]).

The encoder employs dynamic circuits as shown in Figure 3.13. It is fast and relatively compact. The structure is divided into part A and B. The part A reconstructs the address and part B generates the complementary value. The reason to have both polarities is that a single dynamic encoder cannot detect multiple WL errors [47]. The dual encoder can also detect multiple WLS asserted at the same time.

In the layout for the encoder, there are six columns of NMOS transistors for both part A and B. There is no open space between these columns.

A dual dynamic 64 to 6 encoder is used in the scheme. Table 3.2 (after [47]) shows the truth table of the encoder.

Table 3.2 THE TRUTH TABLE OF WL ENCODER PART A AND PART B

WL asserted	Output of part A	Output of part B
WL0	0B000000	0B111111
WL1	0B000001	0B111110
WL2	0B000010	
...
WL63	0B111111	0B000000
No WL is asserted	0B111111	0B111111

3.2.4.1.2 Wren Encoder and NOR Checker

The design of the WREN encoder and NOR checker are similar to that of the WL encoder and NOR checker. No encoding is used and there is no need for complementary circuits owing to the fact that there are very few signals. WREN NOR checker detects errors if any WREN is asserted during a precharge phase. The worst case is when a bad write creeps through when the WL is asserted. Figure 3.14 shows the schematic diagram of the WREN encoder and NOR checker.

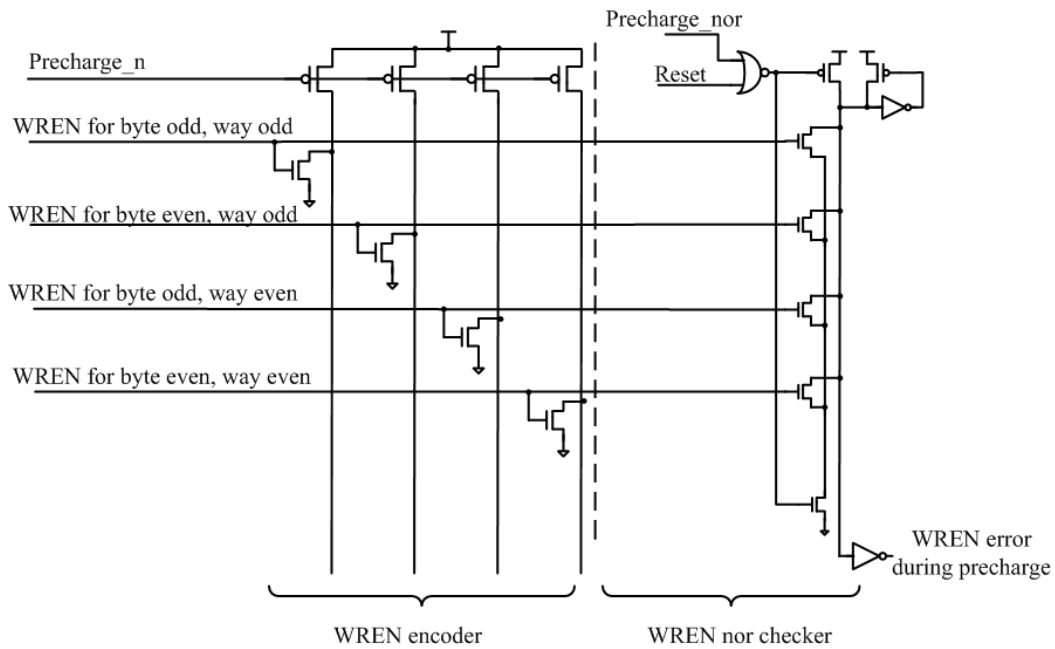


Figure 3.14 The WREN encoder and NOR checker (after [47]).

3.2.4.1.3 BL Precharge Checker

The BL precharge checker as shown in Figure 3.15 monitors the columns connected to even/odd WLs, from the top or bottom sub-bank. When at least one of the four pairs of signals is low during the precharge phase, the checker sets the error flag and a local bit line precharge error is reported. An error detected by BL precharge checker is simulated and illustrated in Figure 3.16.

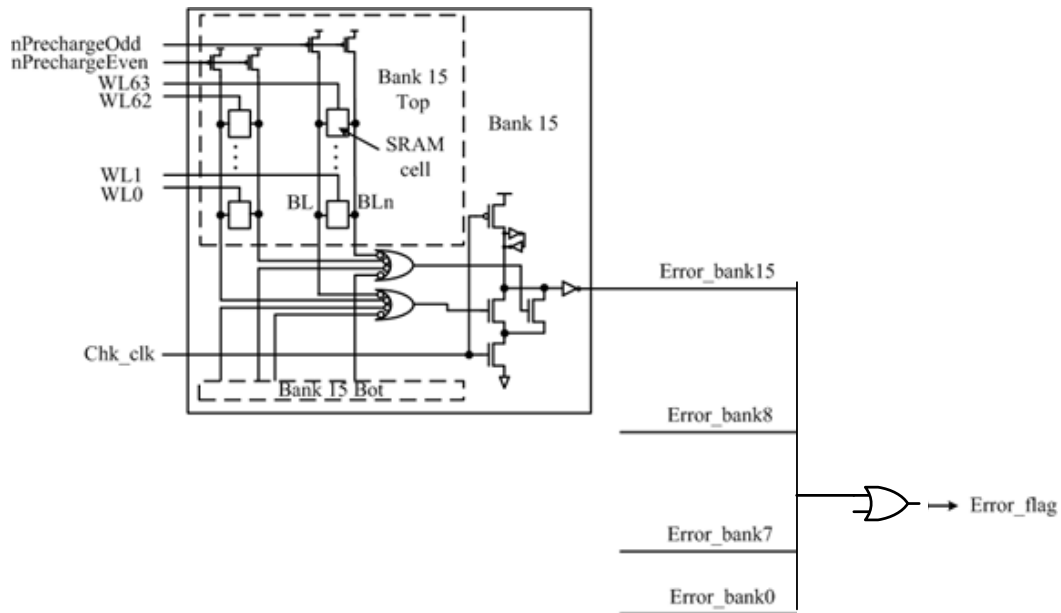


Figure 3.15 BL precharge suppression detection.

In case of an unintended precharge during a read or write, it may cause some contention between the precharge and read or write circuits. If it happens in a write, the write driver is strong enough to overcome the precharge. If it happens in a read, the design ensures that the precharge is stronger and all the outputs are high and thus detectable by the parity checker. Consequently, a checker for unintended BL precharge is not needed.

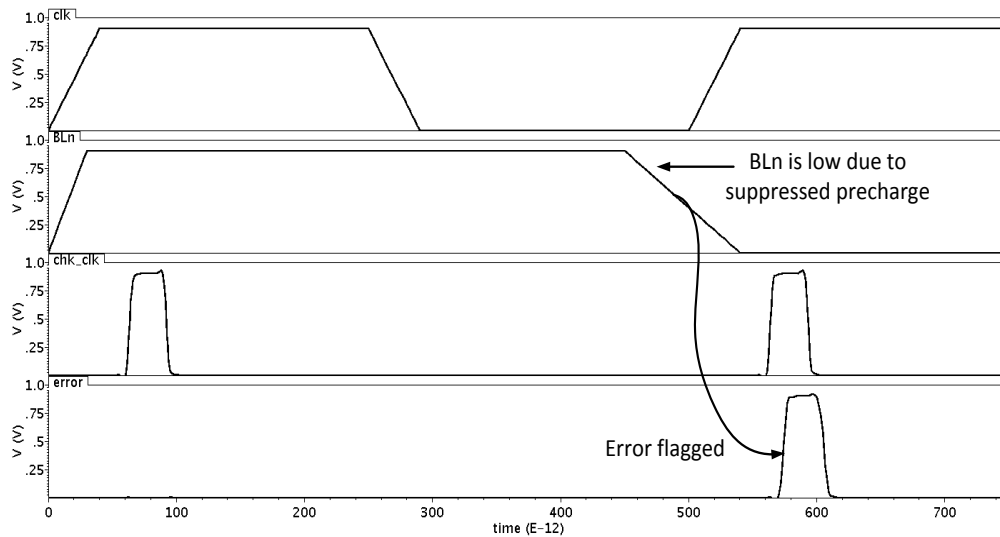


Figure 3.16 Simulation of an error detected on BL precharge.

3.2.4.1.4 BL Read Checker

The BL read checker detects a BL or BLn error caused by a suppressed WL or an asserted precharge signal during a read. A dummy column of special SRAM cells is used and the BL/BLn is monitored. The BL and BLn in the dummy column have almost the same load and precharge driver as the normal columns to duplicate the development of BL and BLn. The value in the SRAM cell is always set to high. During a read, the BL is high and BLn is low. If both are high an error flag is set by the BL read checker.

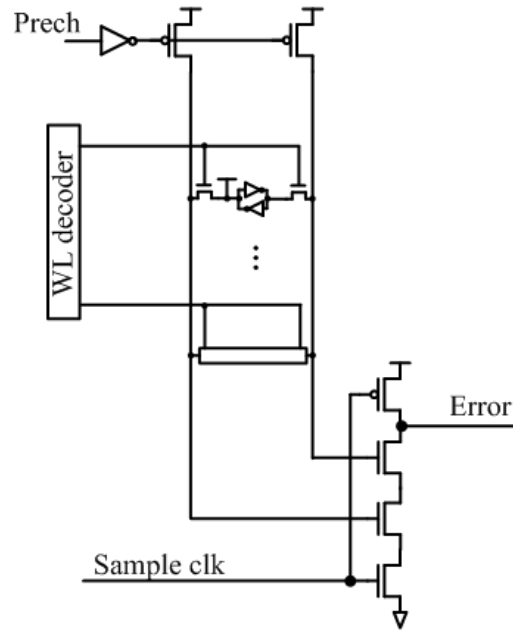


Figure 3.17 Bit line read checker (after [47]).

Multiple SET errors can prevent the BLs from developing completely. For example, the BL precharge can spike or the case when the WL development is delayed. Such an irregularity can cause an incomplete or incorrect read. The read checkers reside in each sub-bank. The error flags from each bank are combined similarly to the BL precharge suppression detection in Figure 3.15. The bit line read checker is shown in Figure 3.17. Its simulated function is illustrated in Figure 3.18.

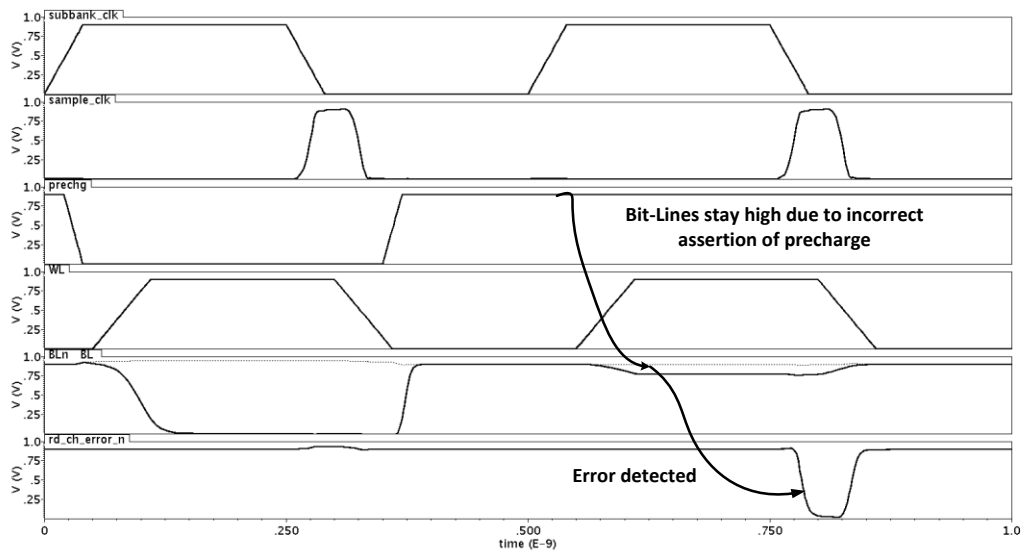


Figure 3.18 Simulation of an error detected on BL during a read.

3.2.4.1.5 Write Checker

The write checker constitutes a column of dummy cells. The cells in the column are connected to a grounded WL. There are 31 dummy cells and one special SRAM cell at the very bottom of the column. In a precharge stage, the special cell is written with logic “0” and during every write stage, logic “1” is written to the special cell. An inside node is monitored for a low state after the write operation. If the node is not low, the write error flag is set. There is a write checker (see Figure 3.19) in each sub-bank. The error flags from all sub-banks are combined in the same way as for the read checker.

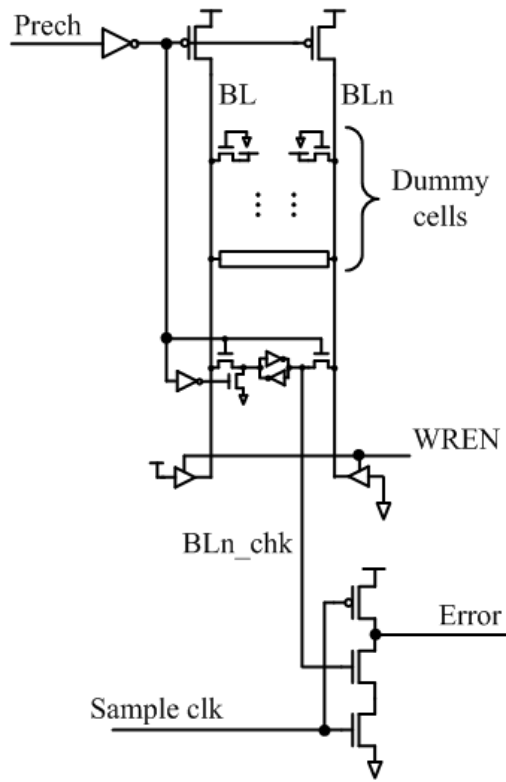


Figure 3.19 Write checker. (After [47])

Figure 3.20 shows simulation waveforms of this checker. A wrong precharge causes both BL and BLn high during a write reporting a write error.

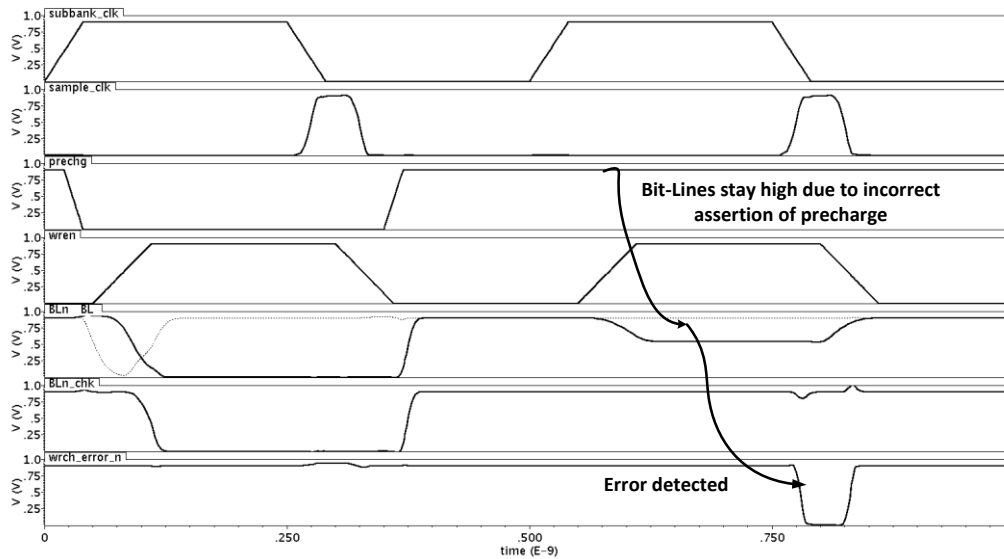


Figure 3.20 Simulation of an error detected during a write.

3.2.4.1.6 SRAM cells used in Read/Write Checkers

The SRAM cells used as dummy for the read and write circuitry are not the standard IBM foundry cells used throughout the design. This was a sort of shortcoming as foundry cells cannot be modified. So the smallest logic cell with the same beta ratio as the IBM foundry cell used was designed and used. Though this may not track exactly the foundry SRAM cells but this is the closest we can get. There is roughly 2 times area overhead from the foundry cell. This is due to the non SRAM DRC rules which are used for the logic cell.

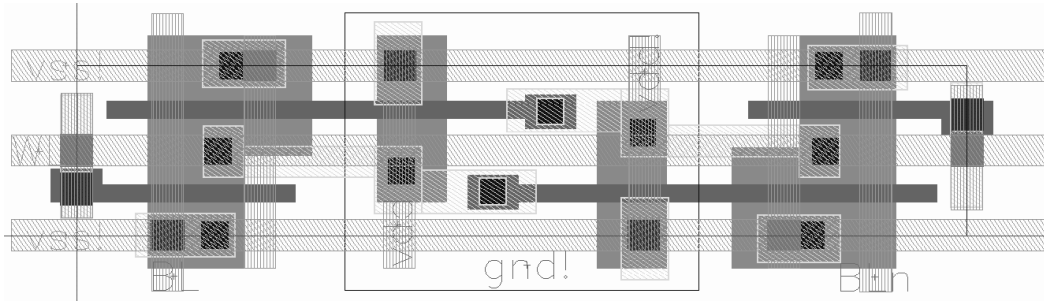


Figure 3.21 Layout of the Logic SRAM cell

Table 3.3 TRANSISTOR SIZES USED FOR THE LOGIC SRAM CELL

PROCESS	IBM 45NM SOI	
DEVICE	WIDTH (μM)	LENGTH (μM)
ACCESS NFET	0.209	0.40
PULL DOWN NFET	0.361	0.40
LOAD PFET	0.152	0.40

Table 3.3 lists the SRAM cell device sizes of all transistors in the 45 nm process.

Table 3.4 READ MARGIN FOR THE DUMMY SRAM CELL FOR 45 NM PROCESS

PROCESS (TYPICAL & WORST CASE (FS))	IBM 45NM SOI	
VDD(V)	TEMP(C)	SNM(V)
0.9	85	0.093
0.72	125	0.054

Table 3.4 gives us the typical and worst case read margins for the custom cell.

Figure 3.23 and Figure 3.22 shows the worst case monte-carlo read and write margins (0.38V) respectively for the logic cell. The worst condition for write is a slow NMOS and a fast PMOS. The worst case for read is a fast NMOS and slow PMOS.

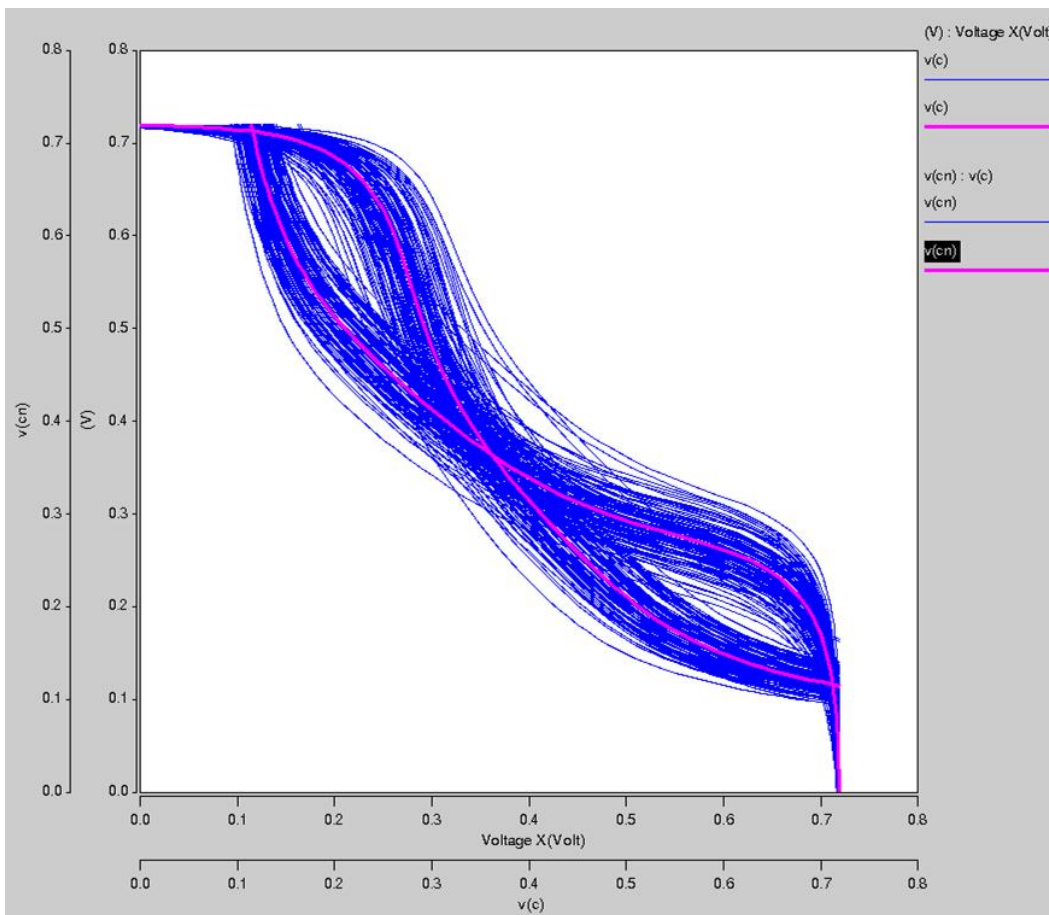


Figure 3.22 Monte Carlo simulations for worst case read SNM analysis

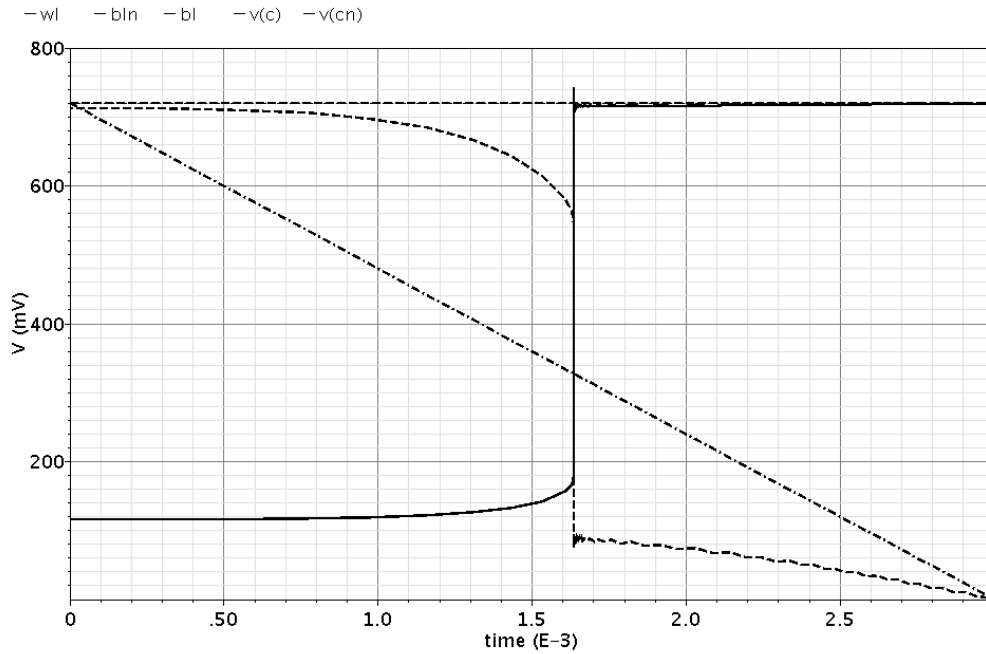


Figure 3.23 Worst case write margin of the Logic SRAM cell

3.3 Tag Array Design Details

The tag is comprised of four ways, with four banks in each way and two sub-banks (top and bottom sub-bank) in each bank. Each way has dual redundant tag comparison and hit generation logic, so that SETs on them can be detected. There are 32 rows and 28 columns in a sub-bank. The components of each way are shown in Figure 3.24. The tag rows use interleaved layout.

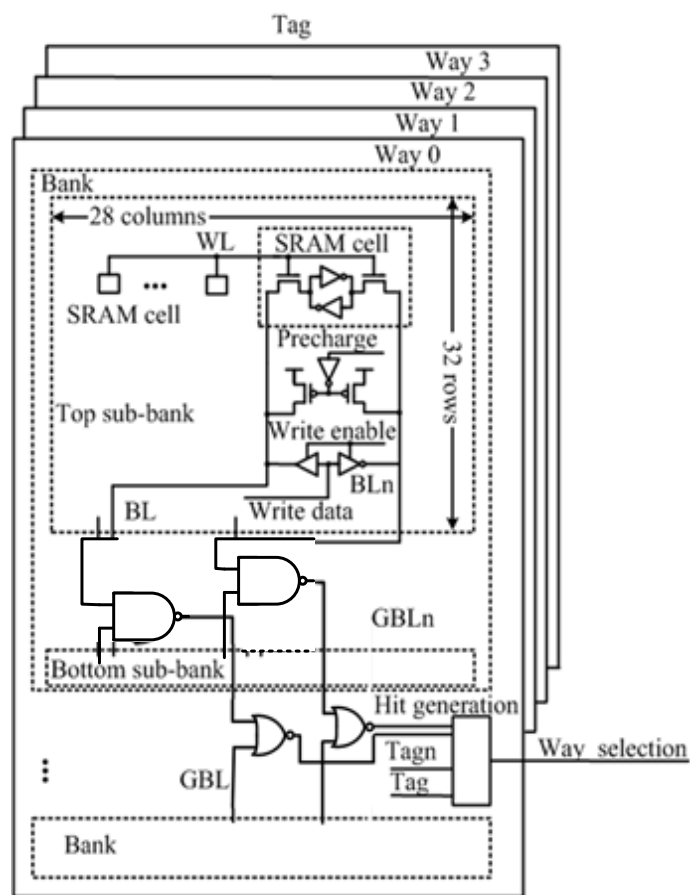


Figure 3.24 Basic diagram of tag.

3.3.1 Tag Array Circuits

The tag array contains the address pertaining to each line in the cache. It should be faster than the data array, since it compares the stored address and selects the appropriate way in the data array. This happens in parallel with the data read. The tag array has WL decoders, SRAM cells columns, precharge for BL, and write circuits which are the same as those used in the data array. In each sub-bank, there are 32 rows and 28 columns of SRAM cells. Unlike the 90nm design the bit lines are not split into 2 having 16 cells each. The bit lines have 32 cells connected to them as shown above. This was done as the 45nm process is on

SOI and has low parasitics compared to a bulk process and there would be minimal penalty on bit line development. This allowed the use of a faster 2 input nand single ended sense (19 ps from bit line discharge to nand sense output) instead of a 4 input nand (22 ps from bit line discharge to nand sense output). This results in lower capacitance on the sense amplifier side allowing for faster speeds and lower power consumption. This simplified the layout and increased array efficiency. The data array still uses the 4-input nand sense owing to the dual redundancy nature of the architecture. The critical path in the cache is the hit generation, so a faster BL development sense is necessary to speed up the way selection signals generation.

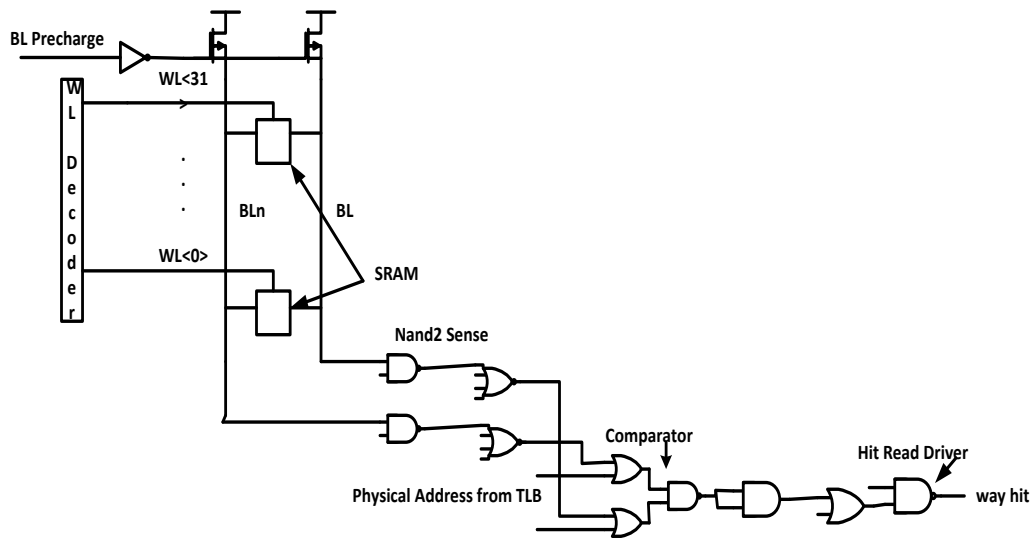


Figure 3.25 Schematic of the tag critical path.

The following sections focus on the circuits that are different from the data array.

3.3.2 Hit Generation

The tag address read from SRAM cells is compared to the physical address in the hit generation circuits, as shown in Figure 3.25. The XNOR gate that takes differential inputs, implemented as a complex CMOS OAI (Or-And-invert) gate, compares one of the 20 bits in the address. The OAI outputs thus generated are ANDed to generate the way hit signal which then goes to the hit read drivers to generate the way selects to the data array.

3.3.3 Error Checking Circuits

The error checking circuits in the tag are similar to those in the data array. The BL precharge, read and write checkers have the same structure as that in the data array. Error flags from the 16 banks are ORed together through static gates and a static NAND gate. When compared to the original 90 nm design, there are 2 pins less in this design as the bit lines have 32 cells instead of 16 cells.

3.3.4 Write and precharge circuitry

The precharge and write enable signals are generated from the subbank clock. The delays are adjusted by sizing the gates in the logic path much like the data array circuits. The logic is shown in Figure 3.26. The non-segregation of the SRAM cells on the bit lines have reduced and simplified the driver circuitry considerably in the tag part.

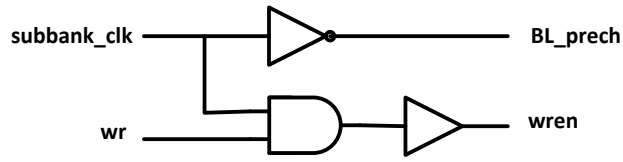


Figure 3.26 Schematic of the precharge and wren circuitry.

3.3.5 Dual Redundant Hit Generation

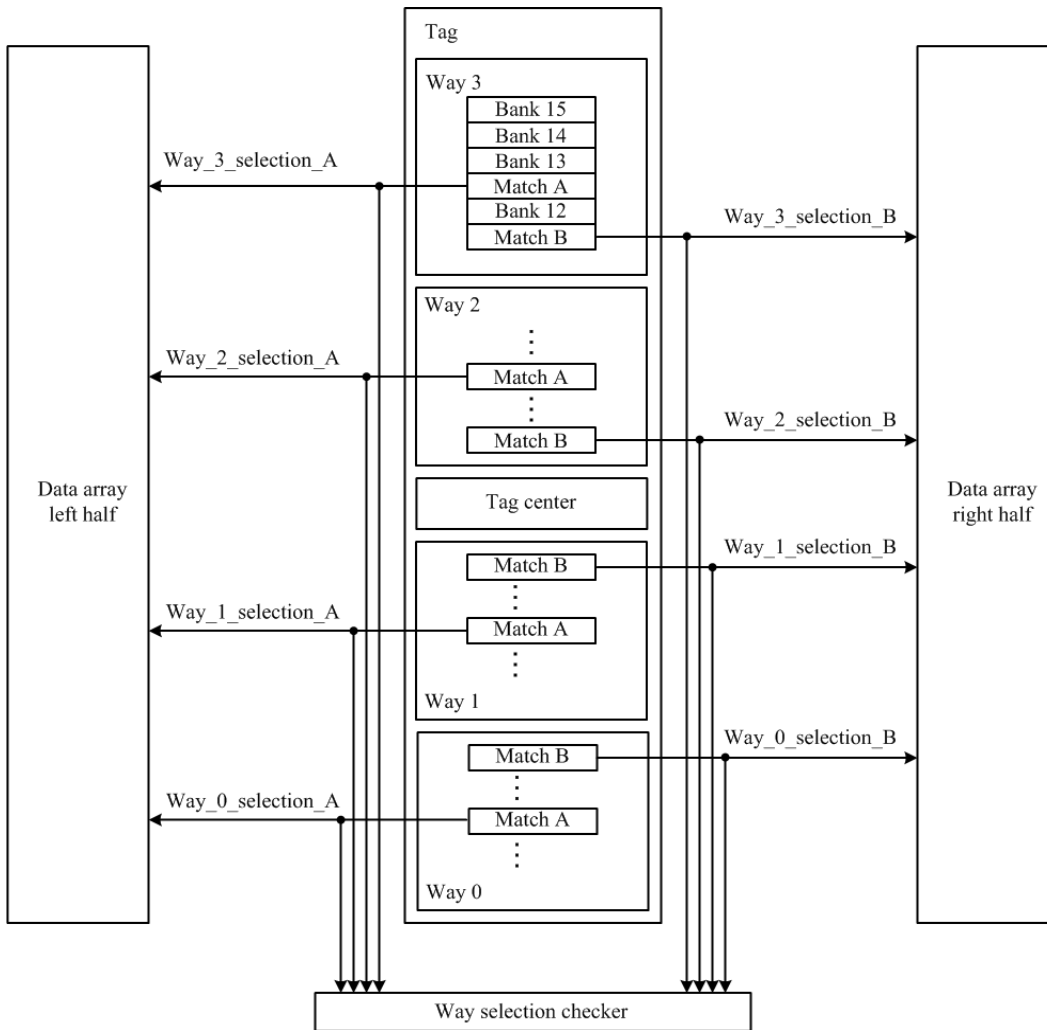


Figure 3.27 Dual match generation of tag (after [47]).

To detect a SEE-induced error in the hit generation circuits, the hit generation logic is designed as dual redundant in each way, named Match A and Match B as illustrated in Figure 3.27. The way select signals from the dual redundant hit generation circuits can be checked against each other. If there is a mismatch, an error is flagged. The checkers for these comparisons lie outside of the whole cache and are not part of the cache design. The dual hit generation circuits are connected to each half of the data array. The load on the way selection circuits is halved than what it would otherwise be. Therefore the drivers of the way selection circuits are sized smaller. The area overhead of the dual hit generation circuits would be less than 30% larger than its single hit generation circuit counterpart, with similar load and performance specifications.

Chapter 4. ARRAY LAYOUT

4.1 Tag Bank Layout

The Tag Bank is the smallest macro in the Tag Array design. It consists of a top sub-bank and a bottom sub-bank. It consists of all the SRAM arrays and checker circuits along with the drivers integrated all together. The tag bank is done full custom layout. The layout snapshot of the tag bank is shown in Figure 4.1. All the checkers as well as the sense amplifiers lie in between the whole bank. There is also sufficient area as and when modifications are desired.

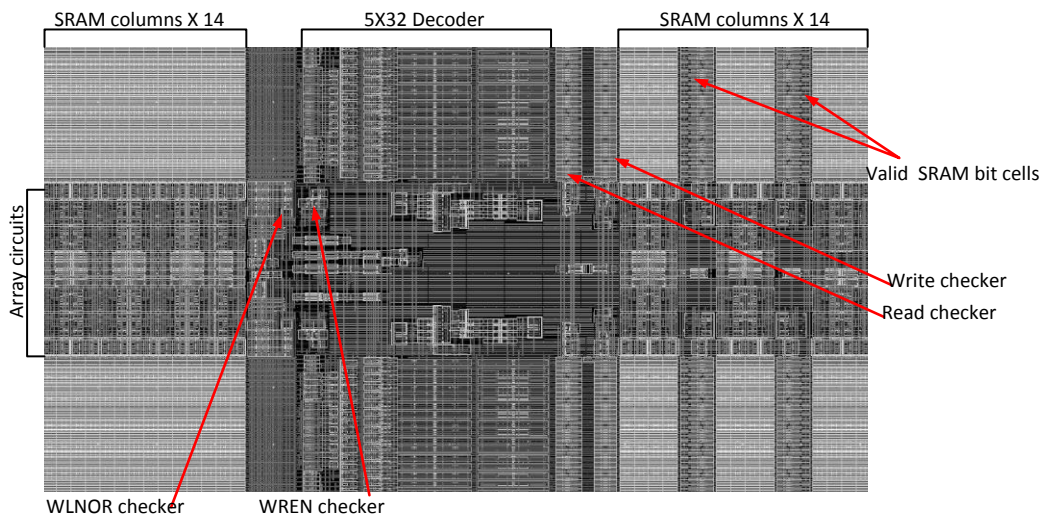


Figure 4.1 Snapshot of the tag bank of the tag array

4.1.1 Layout Interleaving

The tag consists of 28 bits in a row. This includes 20 address bits, four status bits, and four parity bits. They are divided into four parity groups. The distance between bits of the same parity group is the width of three cells or $3.72\mu\text{m}$.

There are two dual redundant valid bits such that they could be protected from SEE. They are cleared independently from other bits in the tag, during power up initialization or global cache invalidation operation. The dual valid bits are output at the same time as and when the cache does a read operation. If either of the valid bits is low, the output is invalid. The distance between the dual valid bits are separated by three SRAM cells. The distance between them is lesser than that was for the previous 90nm design. This is owing to the smaller foundry cells but may be good enough to withstand a particle strike without affecting both the bits.

4.2 Data Bank Layout

The Data Bank like the Tag Bank is the smallest macro in the Data Array design. It also consists of a top and a bottom sub-bank. It consists of all the SRAM arrays and checker circuits along with the drivers integrated all together. The data bank is done full custom layout. The layout snapshot of the data bank is shown in Figure 4.2. All the checkers as well as the sense amplifiers lie in between the whole bank. The data bank SRAMs have dual word lines for interleaved layout between the even and odd data bytes. There is also sufficient area as and when modifications are desired.

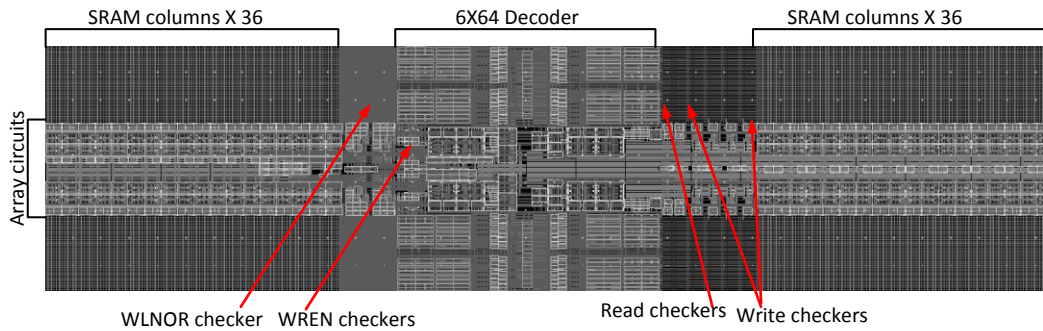


Figure 4.2 Snapshot of the data bank of the data array

4.2.1 Interleaved Layout and Parity Protection

Each row in the left or right half data array contains 72 bits. They are divided into parity groups of eight. There are eight data bits and one parity bit in each group. The eight data bits in the same parity group belong to the same way, WL and byte [47]. Thus bits belonging to different parity groups are interleaved in layout. The distance between bits of the same parity group is same as the width of seven cells or $8.68\mu\text{m}$ wide. The parity checker checks for any single upset within a parity group. In the data array, alternating even and odd parity is computed for consecutive bytes within each word.

4.3 Layout Techniques

One of the most effective layout technique to achieve SEE hardness is to maintain the necessary critical node spacing between cells signals, or blocks that belong to the same parity group or redundant to each other, as emphasized before. The cells in each column are interleaved to maintain this distance. All connections to the cells stay in the same column as the cells. Therefore, all the related circuits also maintain the same critical node distance as the SRAM cells. All redundant

signals, such as the tag hit signals from the tag, maintain the same or greater distance in the layout. The more the spacing, the less susceptible the design is to MBUs. A larger logic SRAM cell is a great advantage in this regard.

4.4 Techniques to Reduce Power

Low power is another goal in the cache design. There are a lot of techniques used widely to reduce power dissipation. One of them is to gate clocks to reduce the activity factor of each sub-bank. Another one is to reduce wire lengths and hence the signaling power dissipation by optimizing the floor plan. We can make use of static CMOS logic wherever possible to reduce power trading a little bit on performance. Also it depends on the micro-architecture of the system. If the architecture is sub-optimal, any number of optimizations would be of no use and would still lead to a poor design or product.

Gated clocks are used to generate each sub-bank clock, as indicated in Figure 4.3. Only banks that are active would receive the clocks. For other sub-banks, the clocks are inactive.

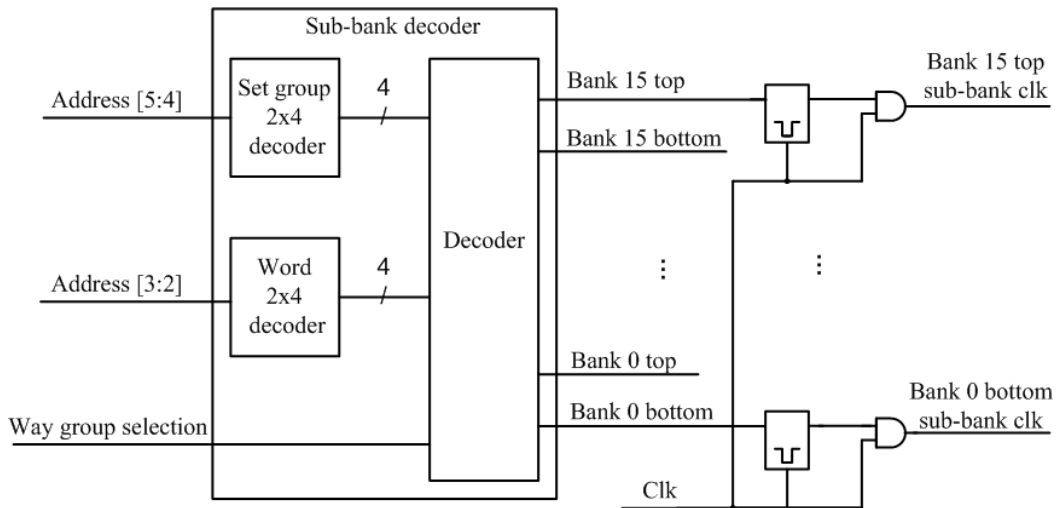


Figure 4.3 Gated clock in the generation of sub-bank clock (after [47]).

Interconnects contribute to the majority of the dynamic power dissipation in the arrays. So reducing the length of interconnects is a necessary part to reduce dynamic power loss. The floor plan for the design is such that the wires are shorter, since the tag arrays are located in the middle of two half parts of the data array, so the hit signals travel a smaller distance. In the data array, the way multiplexers are placed in the center of the array to minimize the delays due to longer wires.

Chapter 5. FULL CACHE LAYOUT

Custom Circuit design has always thrown questions with regards to reliability, resource availability and more importantly the time to design and market such designs. On the other hand ASIC methodologies fall short of the performance mark when compared to custom methodologies while giving better time to market. Achieving custom performance using ASIC techniques is a big challenge for the chip making industry. The blend of the two techniques where the ease of circuit design using standard cells with appropriate sizing and optimal logic implementation, along with the ease of placement and routing of the standard cells, seems to be a captivating solution [53]. This allows better productivity while not compromising much on circuit performance and also leads to lower power and area. One another advantage is the ease of portability from one technology node to another. The full cache is designed using both the methodologies to get the desired performance without much complexity. The proposed methodology used to complete the full cache was actually conceived by Satendra Kumar Maurya as part of his PhD thesis [53].

5.1 The structured methodology flow

One requires a thorough knowledge and analysis of the design architecture to gain maximum custom circuit performance. Figure 5.1 shows the proposed flow. The first step in custom circuit design is to implement the circuit design on a schematic editor using macro and/or standard cells. Special custom standard cells required are also added to the libraries. Design analysis across multiple corners is

subsequently carried out for timing verification using timing tools like Synopsys Prime Time. At this point, the design is not routed, and the wire parasitics are not included. However, the timing tool checks for the proper drive strengths of standard cells by analyzing the transitions at every node in the critical path and gives an optimistic result which gives us an idea of delays post route.

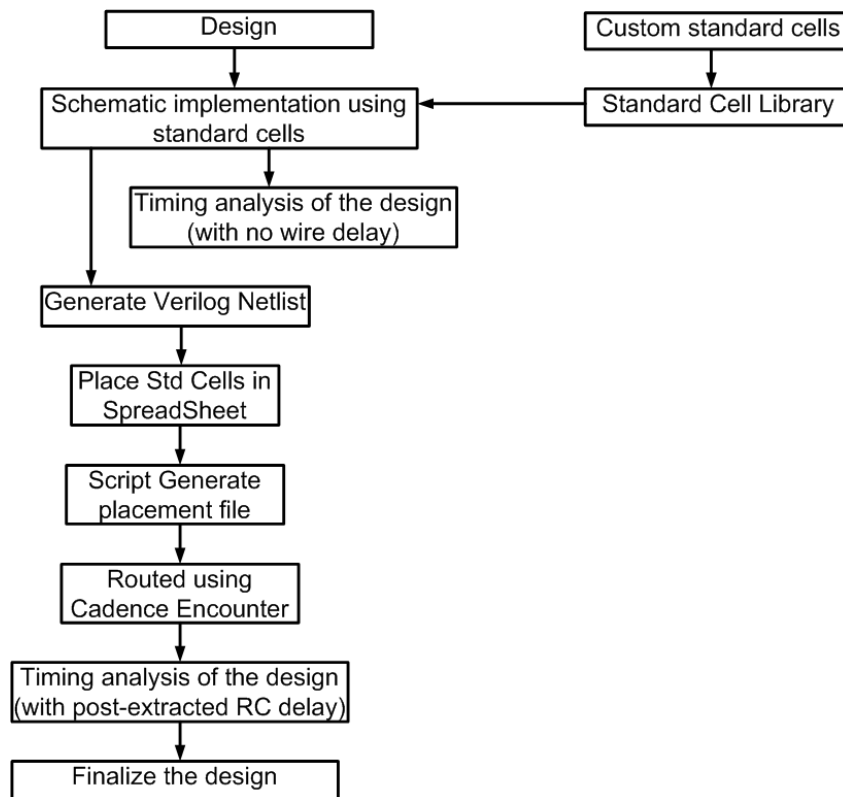


Figure 5.1 Structured methodology flow (after [53]).

Once design at the schematic level is finalized, the floor-planning of these standard cells is carried out. This is done by placing the cells in a spreadsheet, each cell of which corresponds to each instantiation of the standard cell in the design. Since global wires between physical modules are a dominant portion of

the total path delay, proper placement of the standard cells and their drivers make the path delay within the design limit, thus achieving high speed. The spreadsheet is then read by Perl scripts that convert it into a placement file readable by the Cadence Encounter routing tool for actual placement of the standard cells in its floor-plan. The floor-plan in the Cadence Encounter can be analyzed and the standard cells can be replaced to obtain better area. By simply swapping or shifting the placed cells in the spreadsheet and re-running the Perl scripts to obtain the new placements, the task is productive and allows designer flexibility of moving the cells across for area and delay optimization [53].

As soon as the Encounter floor-plan is settled, the routing of all the cells is carried out using the Encounter auto-route tool. This is the step where the ASIC design methodology is used. Since the algorithms for routing are reasonably optimal and regular, they can give similar results as manual routing but at far less effort. The netlist and the timing information (that captures the estimated wire delays) from the routed block can then be extracted. The timing tool then re-analyzes the block but this time with estimated wire RC delays. Changes, if required, can be carried out in the spreadsheet for placement changes or in the schematic for better drive strengths in case the cell is driving long wire. Since these steps are fast and require less effort, the proposed methodology generates a design with better performance, regularity, repeatability, and ease of portability. The main work for the designer is to develop the micro-architecture of the design and circuit design as represented by the schematic.

5.1.1 Steps involved in the structured flow

The steps involved in the flow are mentioned below [53].

1. Create the flattened netlist file using the Verilog-XL simulator in Cadence Virtuoso.
2. Parse all the necessary Library Exchange Format (LEF) files required for the design. They contain physical information for the standard as well as the macro cells. They are generated for custom specific cells using Cadence Abstractgen tool.
3. Plan the layout floorplan using excel sheet and generate the layout information using perl scripts. It contains a list of all cells in their hierarchical as well as location mapping.
4. Finally get the placement information file which can be zipped and loaded into Cadence Encounter tool for placement. It also gives you a rough estimate of the block width and height.

5.1.2 Characterizing macro blocks

To integrate custom design into the structured flow methodology, the tag bank and the data bank which form the smallest macro at the array level were characterized to get the respective .lib files. These library files contain timing information for the macro blocks. The tool used was Synopsys Nanotime which is a transistor level characterization tool. The characterization was performed on the post extracted netlist (extracted using Calibre PEX from IBM). Additionally

HSPICE simulations were performed on the post layout netlist to validate the timings.

5.2 Cache layout

5.2.1 The data array layout

The layout of the data array was carried out in two steps. First, the banks were built independently full custom and then the banks, along with standard cells were put together, to complete the layout. All the 16 banks used in the design are identical. Each of the data banks consists of two sub-banks: top and bottom sub-bank. The custom design achieves a highly compact design. Moreover, there is dynamic logic in the bank, which is best laid out manually. The fully laid out bank was then characterized to generate the library characterization file (*i.e.*, .lib) for the bank. The rest of the cache layout was carried out by placing the banks and the standard cells in the structured flow mentioned before. The banks were treated as macro blocks and their placement was defined in the spreadsheet. Precautions were taken to ensure the RHBD required bit interleaving and maintenance of the spatial distance between traversing signals. Figure 5.2 shows the placement details of the data array. The standard cells are placed in between every two banks. This reduces the wire run for the signals and hence reduces the excessive loading of standard cells. The standard cells are placed in a manner to maintain a laminar flow of the signals to and from the banks. Since the banks maintain the proper bit interleaving and spatial separation between the SEE susceptible signals,

this way of placing the standard cells automatically ensures the radiation hardening criteria.

Once the placement of the standard cells and the banks are finalized in the Encounter floorplan, the design is routed using the Nanoroute feature of Encounter. The design is saved in gdsII format and is then imported back into the Cadence Virtuoso environment. The complete layout of the data array is shown in Figure 5.3. The layout is compact, regular and also since it is generated through the script, any modifications can be applied to the design with reasonable design and characterization effort.

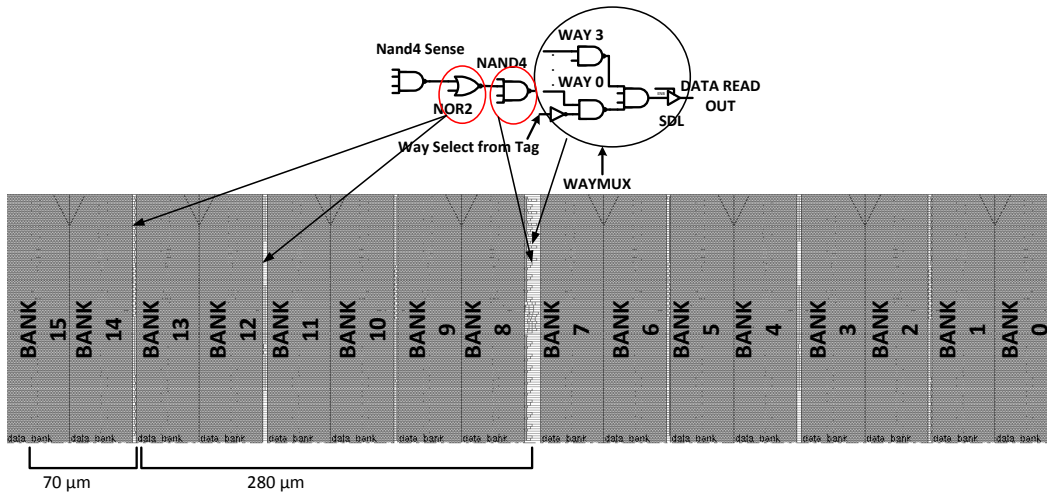


Figure 5.2 Data banks and standard cell placement

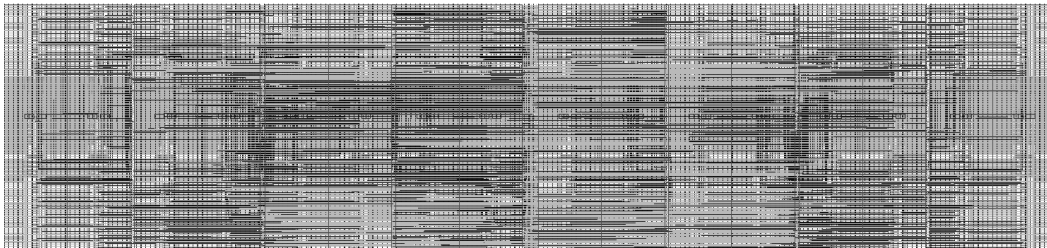


Figure 5.3 Fully routed data array

5.2.2 The tag array layout

The complete layout of the tag array was carried out similar to the data array layout, in two steps. First banks were independently built full custom and then banks along with standard cells were put together to get the complete layout. All the 16 banks used in the design are similar. first these banks were fully custom built. Each of the tag banks consists of two sub-banks: top and bottom sub-banks. The custom design lets to achieve highly compact design and also because there are dynamic logics in the bank which are best laid out manually. Once the layout of the bank is finalized, it was characterized to generate the library characterization file for the bank. Once the banks are ready, the rest of the layout was carried out by placing the banks and the standard cells into the spreadsheet. The banks were treated as macro blocks and their placement was defined in the spreadsheet as per their instance names. Precaution was taken to ensure bit interleaving and spatial distance between the signals traversing from these banks into the standard cells. This spreadsheet is read by the Perl script that generates the Encounter compatible placement file. Figure 5.4 and Figure 5.5 shows the placement details and complete layout of the tag array respectively. The standard cells are placed as strips in a 2-4-2-2-4-2 fashion between the banks. This reduces the wire run for the signals and hence reduces the excessive loading of standard cells. The standard cells are placed in a manner to maintain a laminar flow of the signals to and from the banks. Since the banks maintain the proper bit interleaving

and spatial separation between the SEE susceptible signals, this way of placing the standard cells automatically ensure the radiation hardening criteria.

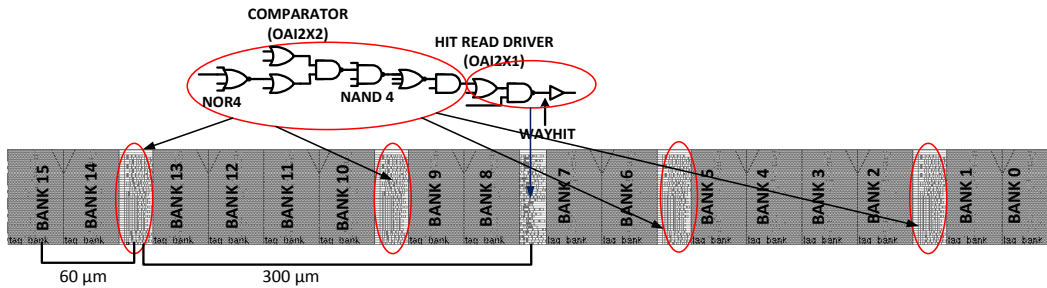


Figure 5.4 Tag banks and standard cell placement

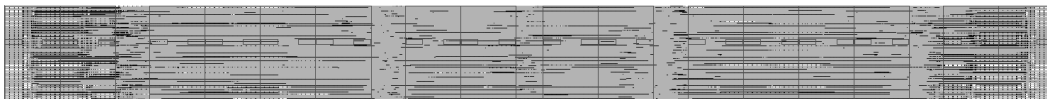


Figure 5.5 Fully routed tag array

5.2.3 Full cache layout

Once the tag and data arrays are fully laid out, Synopsys Primitime is run on the extracted netlist from Encounter to characterize the arrays and to get the .lib file which has all the timing information. To build the entire cache the structured methodology described before is again applied. There will be 1 placement for the tag array and 2 placements for the data arrays which reflect the right and left half data arrays. The floorplan of the entire cache is shown in Figure 5.6. The LEFs for both the completed tag and data arrays are extracted using Abstractgen. The placement information is again generated using scripts and fed into Encounter. Nanoroute is used for routing the entire cache. Once routed, the entire design is extracted and fed as an input to Primitime for timing analysis. The Primitime results are matched to the HSPICE simulation results.

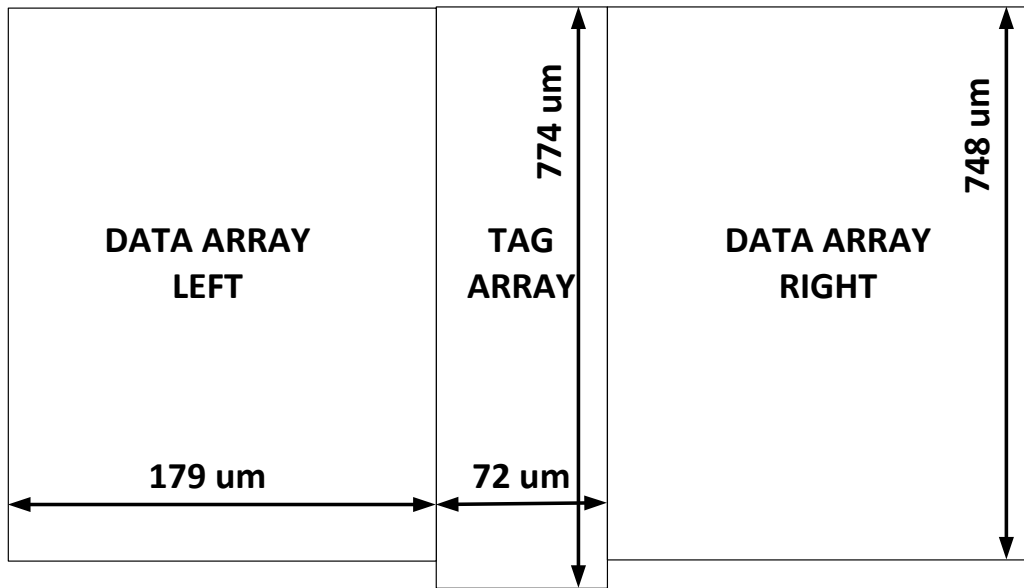


Figure 5.6 Floorplan of the cache

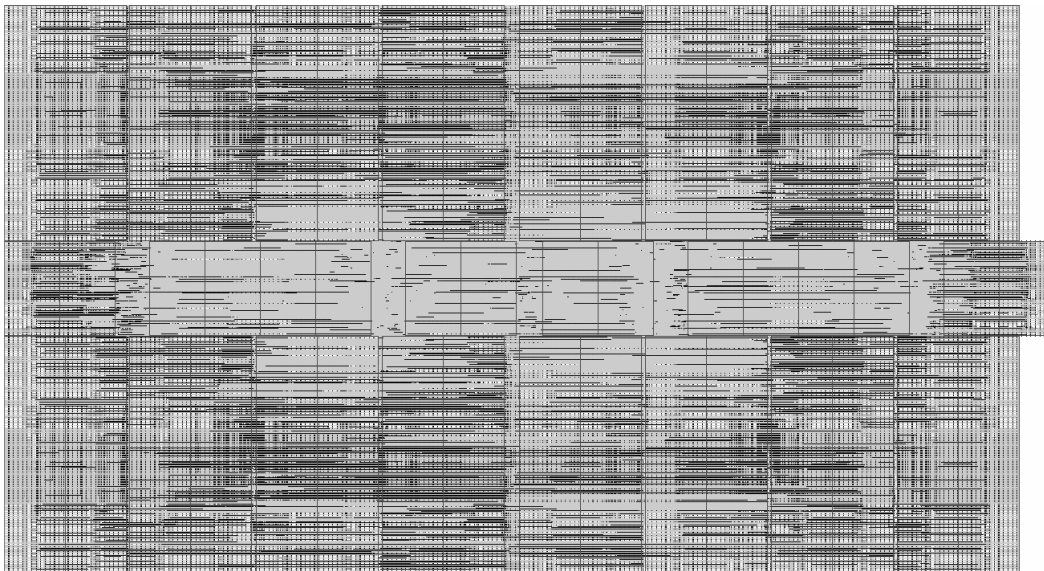


Figure 5.7 The routed entire Cache

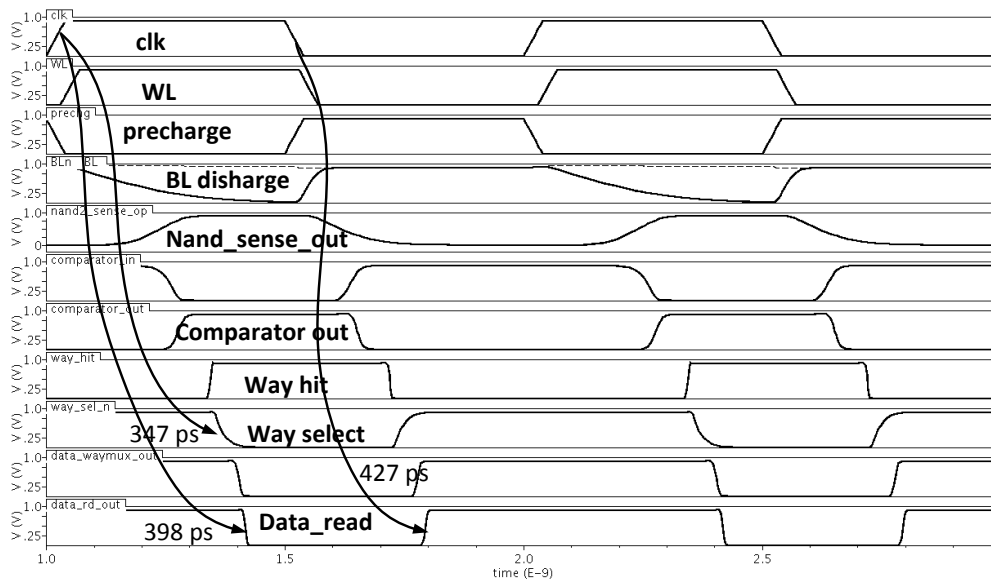


Figure 5.8 HSPICE simulation showing the Cache critical path

Figure 5.7 shows the complete layout of the whole cache. Figure 5.8 gives a simulation waveform snapshot of the critical path of the Cache design. The worst case path delay is around 398 ps guaranteeing that the cache design can work above 2 GHz. Table 5.1, Table 5.2 and Table 5.3 show the delay and area, average power and energy consumed respectively for the 45 nm and 90 nm cache designs.

Table 5.1 DELAY AND AREA FOR VARIOUS CACHE DESIGNS

Cache Design	Delay (ns)	Area (mm ²)
90 nm dynamic	0.930	3.422
90 nm static	1.017	3.163
45 nm static	0.398	0.331

Table 5.2 AVERAGE POWER

	45 nm Cache (mW)		90 nm Ccahe (mW)	
	Read	Write	Read	Write
Half data array	14	16	53	49
Tag array	9.6	10	21	19
Cache	37.6	42	127	117

Table 5.3 ENERGY PER OPERATION

45 nm Cache	Energy per operation (pJ)	
	Read	Write
Half data array	14	31.6
Tag array	10.6	20
Cache	38.6	83.2

Chapter 6. CONCLUSION

This thesis provides details about the design and implementation of a radiation hardened by design cache. The main goal is to design a radiation hardened cache which can boast of low power and high performance while keeping portability in mind. Portability refers to the time taken to port a design from one technology node to another in as less time as possible.

RHBD techniques are used to mitigate errors due to SEE. Special checker circuits, interleaved layouts and dual redundancy approaches are employed. The use of IBM foundry SRAM cells leads to better array efficiency and low power when compared to the use of specially made logic SRAM cells which would be very large. Although in the perspective of TID, the smaller size would be a disadvantage and make the design less hard. Special logic SRAM cells were designed for Tag parity bits and other dummy SRAM cells used in the read and write checker circuits. This was mainly owing to the inability to modify the foundry SRAM cells. Dual WLs (wordlines) could be employed on the foundry cell for the purpose of interleaving even and odd data bytes on the foundry SRAM cells. This was especially hard due to the small size of the SRAM cells. The operating speeds can touch above 2GHz which is way higher than any of the commercial RHBD processors available commercially. This was possible due to the optimal use of both dynamic as well as static CMOS circuits. The cache achieves low power by making use of gated clocks to the subbanks. This means that at one time only one subbank would be active and power would be saved as

the rest of the 15 banks would be turned off. Also the placement of the tag array and the 2 halves of the data array on each side of the tag array minimizes the interconnect lengths and standard cells with lesser drives and thus lesser capacitance can be used. This actually reduces the distance between the way hit logic in the tag to the way-muxes in the data array. This also minimizes power. This cache design also makes use of single ended nand sense instead of the analog sense amplifiers used in traditional caches. It also does not make use of column multiplexers. This use of CMOS sense reduces mismatch due to variability and simplifies circuit design. The SOI technology allowed the use of a 2 input nand sense instead of the 4 input nand in the 90nm design. SOI technology has less transistor parasitics than bulk and allows for faster design and thus less parasitics on the bitlines leads to faster bitline development speed and reduced complexity of peripheral circuits. The original design had partitioned the bitlines into top and bottom bitlines to improve bitline development speed.

The use of the structured ASIC methodology flow to build high performance arrays has also been exploited in this thesis. The performance of custom circuits along with the flexibility and time to market capability of the ASIC structured flow leads to the development of robust designs. The flow gives the designer the freedom to add more custom cells to the library as and when need arises. The data and tag banks as macros can be characterized and their timing data can be used in this flow. They are like standard .lib files which can be used in Encounter or Primitime for layout or timing analysis purposes. In a traditional

ASIC flow you would not be able to use special cells as you are restricted to only a certain standard cell library. The functional verification is done in a RTL model vs schematic model approach. Here the vectors from a behavioral model written in HDL is used to test the post layout extracted netlist and checked for functional bugs. Cadence Ultrasim tool was used for this purpose. Once the functional verification is done the layout is done using the spreadsheet format where the macros and standard cells are placed for optimum placement and cells in critical paths are placed such as to reduce interconnect parasitics. Once the placement is finalized and reiterated to get the best placement in Encounter, routing is done and the parasitics are extracted as a Primetime readable format. Finally static timing analysis is performed for violations and to determine the delay on the critical path. Any problems on the path like longer delays or low drives can be fixed by sizing the standard cells and reiterating the entire flow till optimum results are achieved. This can be done in considerably short amount of time unlike for a full custom design. Finally many novel design techniques were involved. The judicious use of dynamic logic and SDLs can make the design highly optimized. The uses of SDLs were an integral part of the design. The ported 45 nm design gives significant savings in terms of area and power as well as reduction in overall area when compared to the 90 nm design.

REFERENCES

- [1] R. Koga, S. D. Pinkerton, S. C. Moss, D. C. Mayer, S. Lalumondiere, S. J. Hansel, K. B. Crawford, and W. R. Crain, "Observation of single event upsets in analog microcircuits," *IEEE Trans. Nucl. Sci.*, vol. 40, no. 6, pp. 1838–1844, Dec. 1993.
- [2] R. Ecoffet, S. Duzellier, P. Tastet, C. Aicardi, and M. Labrunee, "Observation of heavy ion induced transients in linear circuits," *Proc. IEEE NSREC Radiation Effects Data Workshop Record*, pp. 72–77, 1994.
- [3] R. Harboe-Sorensen, F. X. Guerre, H. Constans, J. Van Dooren, G. Berger, and W. Hajdas, "Single event transient characterization of analog ICs for ESA's satellites," in *Proc. IEEE 5th Eur. Conf. Radiation and Its Effects on Components and Systems*, 2000, pp. 573–581.
- [4] B. E. Pritchard, G. M. Swift, and A. H. Johnston, "Radiation effects predicted, observed, and compared for spacecraft systems," *Proc. IEEE NSREC Radiation Effects Data Workshop Record*, pp. 7–17, 2002.
- [5] Narayanan. V, Xie. Y, "Reliability concerns in embedded system designs," *Computer*, vol.39, no.1, pp. 118- 120, Jan. 2006.
- [6] Hughes, H.L, Benedetto, J.M, "Radiation effects and hardening of MOS technology: devices and circuits," *Nuclear Science, IEEE Transactions on*, vol.50, no.3, pp. 500- 521, June 2003.
- [7] Barth, J.L.; Dyer, C.S.; Stassinopoulos, E.G.; , "Space, atmospheric, and terrestrial radiation environments," *Nuclear Science, IEEE Transactions on* , vol.50, no.3, pp. 466- 482, June 2003.
- [8] Stassinopoulos, E.G., Raymond, J.P., "The space radiation environment for electronics," *Proceedings of the IEEE*, vol.76, no.11, pp.1423-1442, Nov 1988.
- [9] Gussenhoven, M.S.; Mullen, E.G.; Brautigam, D.H., "Improved understanding of the Earth's radiation belts from the CRRES satellite," *Nuclear Science, IEEE Transactions on*, vol.43, no.2, pp.353-368, Apr 1996.
- [10] H. L. Hughes and R. R. Giroux, "Space radiation affects MOSFET's," *Electronics*, vol. 37, p. 58, 1964.

- [11] A. R. Frederickson, "Upsets related to spacecraft charging," IEEE Trans. Nucl. Sci., vol. 43, no. 2, pp. 426-441, 1996.
- [12] Kerns, S.E.; Shafer, B.D.; Rockett, L.R., Jr.; Pridmore, J.S.; Berndt, D.F.; van Vonno, N.; Barber, F.E.; , "The design of radiation-hardened ICs for space: a compendium of approaches," Proceedings of the IEEE , vol.76, no.11, pp.1470-1509, Nov 1988.
- [13] Hsieh, C. M.; Murley, P. C.; O'Brien, R. R., "Dynamics of Charge Collection from Alpha-Particle Tracks in Integrated Circuits," Reliability Physics Symposium, 1981. 19th Annual, vol., no., pp.38-42, April 1981.
- [14] Mavis, D.G.; Eaton, P.H., "Soft error rate mitigation techniques for modern microcircuits," Reliability Physics Symposium Proceedings, 2002. 40th Annual , vol., no., pp. 216- 225, 2002
- [15] Barnaby, H. J.; , "Total-Ionizing-Dose Effects in Modern CMOS Technologies," Nuclear Science, IEEE Transactions on , vol.53, no.6, pp.3103-3121, Dec. 2006
- [16] Saggese, G.P.; Wang, N.J.; Kalbarczyk, Z.T.; Patel, S.J.; Iyer, R.K.; , "An experimental study of soft errors in microprocessors," Micro, IEEE , vol.25, no.6, pp. 30- 39, Nov.-Dec. 2005
- [17] Dodd, P.E.; Massengill, L.W.; , "Basic mechanisms and modeling of single-event upset in digital microelectronics," Nuclear Science, IEEE Transactions on , vol.50, no.3, pp. 583- 602, June 2003
- [18] Hsieh, C.M.; Murley, P.C.; O'Brien, R.R., "A field-funneling effect on the collection of alpha-particle-generated carriers in silicon devices," Electron Device Letters, IEEE, vol.2, no.4, pp.103-105, April 1981.
- [19] Karnik, T.; Hazucha, P.; , "Characterization of soft errors caused by single event upsets in CMOS processes," Dependable and Secure Computing, IEEE Transactions on , vol.1, no.2, pp. 128- 143, April-June 2004
- [20] S. J. Heileman, W. R. Eisenstadt, R. M. Fox, R. S. Wagner, N. Bordes, and J. M. Bradley, "CMOS VLSI single event transient characterization," IEEE Trans. Nucl. Sci., vol. 36, no. 6, pp. 2287-2291, 1989.
- [21] Benedetto, J.; Eaton, P.; Avery, K.; Mavis, D.; Gadlage, M.; Turflinger, T.; Dodd, P.E.; Vizkelethy, G.; , "Heavy ion-induced digital single-event

- transients in deep submicron Processes," Nuclear Science, IEEE Transactions on , vol.51, no.6, pp. 3480- 3485, Dec. 2004
- [22] Gadlage, M.J.; Schrimpf, R.D.; Benedetto, J.M.; Eaton, P.H.; Mavis, D.G.; Sibley, M.; Avery, K.; Turflinger, T.L.; , "Single event transient pulse widths in digital microcircuits," *Nuclear Science, IEEE Transactions on* , vol.51, no.6, pp. 3285- 3290, Dec. 2004
- [23] Axness, C. L.; Weaver, H. T.; Fu, J. S.; Koga, R.; Kolasinski, W. A.; , "Mechanisms Leading to Single Event Upset," Nuclear Science, IEEE Transactions on , vol.33, no.6, pp.1577-1580, Dec. 1986.
- [24] Musseau, O.; Gardic, F.; Roche, P.; Corbiere, T.; Reed, R.A.; Buchner, S.; McDonald, P.; Melinger, J.; Tran, L.; Campbell, A.B.; , "Analysis of multiple bit upsets (MBU) in CMOS SRAM," Nuclear Science, IEEE Transactions on , vol.43, no.6, pp.2879-2888, Dec 1996.
- [25] Koga, R.; Penzin, S.H.; Crawford, K.B.; Crain, W.R., "Single event functional interrupt (SEFI) sensitivity in microcircuits," RADECS 97. Fourth European Conference on, vol., no., pp.311-318, 15-19 Sep 1997.
- [26] H. J. Barnaby, M. L. McLain, I. S. Esqueda, and X. J. Chen, "Modeling ionizing radiation effects in solid state materials and CMOS devices," IEEE Trans. Circuits Syst. I, vol. 56, pp. 1870–1883, Aug. 2009.
- [27] Lacoce, R.C.; Osborn, J.V.; Koga, R.; Brown, S.; Mayer, D.C.; , "Application of hardness-by-design methodology to radiation-tolerant ASIC technologies," Nuclear Science, IEEE Transactions on , vol.47, no.6, pp.2334-2341, Dec 2000
- [28] P. Dodd and F. Sexton, "Critical charge concepts for CMOS SRAMs," IEEE Trans. Nucl. Sci., vol. 42, no. 6, pp. 1764–1771, Dec. 1995.
- [29] Hindman, N.D.; Clark, L.T.; Patterson, D.W.; Holbert, K.E.; , "Fully Automated, Testable Design of Fine-Grained Triple Mode Redundant Logic," Nuclear Science, IEEE Transactions on , vol.58, no.6, pp.3046-3052, Dec. 2011
- [30] Hindman, N.D.; Pettit, D.E.; Patterson, D.W.; Nielsen, K.E.; Xiaoyin Yao; Holbert, K.E.; Clark, L.T.; , "High speed redundant self-correcting circuits for radiation hardened by design logic," Radiation and Its Effects on Components and Systems (RADECS), 2009 European Conference on , vol., no., pp.465-472, 14-18 Sept. 2009

- [31] C. Weaver, J. Emer, S. Mukherjee, and S. Reinhardt, "Techniques to reduce the soft error rate of a high-performance microprocessor," Proc. ISCA, 2004, pp. 264-27
- [32] Fujiwara, E.; Pradhan, D.K., "Error-control coding in computers," Computer, vol.23, no.7, pp.63-72, July 1990.
- [33] Chen, C. L.; Hsiao, M. Y.; , "Error-Correcting Codes for Semiconductor Memory Applications: A State-of-the-Art Review," IBM Journal of Research and Development , vol.28, no.2, pp.124-134, March 1984
- [34] Yao. X; Clark, L.T.; Chellappa, S.; Holbert, K.E.; Hindman, N.D., "Design and Experimental Validation of Radiation Hardened by Design SRAM Cells," Nuclear Science, IEEE Transactions, vol.57, no.1, pp.258-265, Feb. 2010.
- [35] Knudsen, J. E.; Clark, L. T., "An Area and Power Efficient Radiation Hardened by Design Flip-Flop," Nuclear Science, IEEE Transactions on, vol.53, no.6, pp.3392-3399, Dec. 2006.
- [36] LaBel, K.A.; Gates, M.M., "Single-event-effect mitigation from a system perspective," Nuclear Science, IEEE Transactions on, vol.43, no.2, pp.654-660, Apr 1996.
- [37] R. Lacoce, "CMOS scaling, design principles and hardening-by-design methodologies," presented at the IEEE Nuclear and Space Radiation Effects Conf., Monterey, CA, Jul. 2003, short course.
- [38] D. R. Alexander, "Design issues for radiation tolerant microcircuits for space", Short Course Nuclear and Space Radiation Effects Conf., 1996.
- [39] W. Beauvais, P. McNulty, W. A. Kader, and R. Reed, "SEU parameters and proton-induced upsets," in Proc. Sec. European Conf. on Radiation and its Effects on Components and Systems, pp. 540-545, Sept. 1993.
- [40] W. Massengill, M. Alles, and S. Kerns, "SEU error rates in advanced digital CMOS," in Proc. Sec. European Conf. on Radiation and its Effects on Components and Systems, pp. 546-553, Sept. 1993.
- [41] D. Kobayashi, T. Makino, and K. Hirose, "Analytical expression for temporal width characterization of radiation-induced pulse noises in SOI CMOS logic gates," Proc. IRPS, pp. 165-169, 2009.

- [42] X. Yao, L. Clark, D. Patterson, K. Holbert, "A 90 nm bulk CMOS radiation hardened by design cache memory," *IEEE Trans. Nucl. Sci.*, vol. 57, no. 4, pp. 2089-2097, Aug. 2010.
- [43] T. Calin, M. Nicolaidis and R. Velazco, "Upset hardened memory design for submicron CMOS technology," *IEEE Trans. Nucl. Sci.*, vol. 43, no. 6, pp. 2874-2878, Dec. 1996.
- [44] A. Chandrakasan, W. J. Bowhill, and F. Fox, *Design of High-Performance Microprocessor Circuits*, 2001 :IEEE Press
- [45] R. Bauman, "The impact of technology scaling on soft error rate performance and limits to the efficacy of error correction," *International Electron Devices Meeting*, pp. 329-332, 2002.
- [46] K. Mohr and L. Clark, "A radiation hardened by design register file with lightweight error detection and correction," *IEEE Trans. Nucl. Sci.*, Vol. 54, pp. 1335-1342, August 2007.
- [47] X. Yao, "Radiation hardened high performance microprocessor cache design", Phd dissertation, Arizona State University, 2009
- [48] D.A. Patterson, John. L. Hennessy, *Computer Architecture: A quantitative Approach*. 2nd ed., Morgan Kaufmann Publishers, Inc., 1996, pp. 375-390.
- [49] J. Handy, *The Cache Memory Book*, San Diego: Academic Press, Inc, 1998.
- [50] J.D. Gee, M.D. Hill, D.N. Pnevmatikatos, and A.J. Smith, "Cache performance of the SPEC92 benchmark suite," *IEEE Micro*, vol. 13, issue 4, pp. 17-27, August 1993.
- [51] Intel Corp., Microprocessor quick reference guide, [Online]. Available: <http://www.intel.com/pressroom/kits/quickreffam.htm#core2>
- [52] N. Weste and D. Harris, "CMOS VLSI Design: A circuit and systems perspective", 4th Ed., Addison-Wesley publication.
- [53] S.K. Maurya, "A Structured Design Methodology for High performance VLSI Arrays", PhD dissertation, Arizona State University, 2012

- [54] J. Benedetto, H. Boesch, and F. McLean, "Dose and energy dependence of interface trap formation in Cobalt-60 and X-ray environments," *IEEE Trans. Nuc. Sci.*, vol. 35, pp. 1260-1264, December 1988.
- [55] G. C. Stierhoff and A. G. Davis. "A history of the IBM Systems Journal," *IEEE Annals of the History of Computing*, vol. 20, pp. 29-35, January 1998.
- [56] S. Buchner, P. Marshall, S. Kniffin, K. LaBel "Proton Testing Guidelines," NASA, 2002
- [57] F. B. McLean and T. R. Oldham, "Basic mechanisms of radiation effects in electronic materials and devices," *Harry Diamond Lab. Tech.Rep.*, vol. HDL-TR, pp. 2129, 1987.
- [58] C.S. Dyer, "Space Radiation Environment Dosimetry," IEEE NSREC, 1998, Section II, pp.16
- [59] G.K. Lum, "Hardness assurance for space systems," IEEE NSREC, 2004, Section I, pp.17
- [60] X. Yao, N. Hindman, L. Clark, K. Holbert, D. Alexander and W. Shedd, "The impact of total ionizing dose on unhardened SRAM cell margins," *IEEE Trans. Nuc. Sci.*, vol. 55, pp. 3280-3287, December 2008.