

Thermal Neutrons: a New Threat for Supercomputers and Safety Critical Applications

ABSTRACT

The high performance, high efficiency, and low cost of Commercial Off-The-Shelf (COTS) devices make them attractive for applications with strict reliability constraints. Today, COTS devices are adopted in HPC and safety-critical applications such as autonomous driving. Unfortunately, the cheap natural Boron widely used in COTS chip manufacturing process makes them highly susceptible to thermal (low energy) neutrons.

In this paper, we demonstrate that thermal neutrons are a significant threat to COTS device reliability. For our study, we consider two DDR memories, an AMD APU, three NVIDIA GPUs, an Intel accelerator, and an FPGA executing a relevant set of algorithms. We consider different scenarios that impact the thermal neutron flux such as weather, concrete walls and floors, and HPC liquid cooling systems. Correlating beam experiments and neutron detector data, we show that thermal neutrons FIT rate¹ could be comparable or even higher than the high energy neutron FIT rate.

1. INTRODUCTION

Reliability is one of the most important considerations in the field of High Performance Computing (HPC) [1, 2, 3]. An unreliable system can negatively affect not only the throughput of a computer but also the correctness of operations. Reliability can be increased through redundancies in chip architectures, improved manufacturing processes, transistor layout changes, or other hardening solutions [4]. However this added reliability comes at an added cost in terms of additional engineering, more expensive manufacturing, and added power consumption. This creates a trade off between lower cost and higher reliability such that only specialized safety critical industries, such as aerospace or medical, are willing to pay the additional cost of highly reliable parts. This is in contrast to Commercial Off-The-Shelf (COTS) devices which are generally not built to the highest achievable levels of reliability due to the low margins of the markets that consume these parts. Most consumers of COTS parts are primarily interested in performance and low price. They are typically willing to suffer lower reliability in exchange [4, 5]. The major-

ity of the HPC community builds systems out of COTS parts and there is a constant struggle between the drive for ever increasing compute power and the potential of lower scientific productivity due to lower reliability [6].

In this paper we compare the reliability risk to HPC systems from *high energy* neutrons to that of boron-10 (^{10}B), which makes devices vulnerable to *thermal neutrons* generated from either fast neutrons that have lost energy through multiple interactions [5, 7] or are emitted from naturally occurring radioactive isotopes. ^{10}B has a relatively large capture cross section for thermal neutrons and the resulting excited state of ^{10}B quickly decays into Lithium-7 and a 1.47 MeV alpha particle. It is this high energy alpha particle that is known to contribute to upsets in semiconductors. Eliminating boron all-together or using depleted ^{11}B would make the device immune to thermal neutrons. However, depleted boron is expensive and boron is necessary for the manufacture of modern semi-conductors, so many COTS devices contain ^{10}B . Modern data centers contain large masses of materials that can potentially increase the flux of thermal neutrons, in the form of concrete slab floors, cinder block walls, and water cooling units. In order to accurately estimate the effects of thermal neutrons we deployed a neutron detector to measure the natural background rate variation due to materials used in a modern data center. Our initial measurements indicate that these materials can increase the thermal neutron counts, and thus the COTS device's error rate, by as much as 20%.

The details of how ^{10}B is used in modern chips is proprietary and not publicly available. The only way to evaluate boron concentration in a chip, and the associated increased sensitivity to thermal neutrons, is through controlled radiation exposure. We studied the effects of fast and thermal neutrons on DDR3 and DDR4 memories, an AMD Accelerated Processing Unit (APU), three NVIDIA GPUs, an Intel accelerator, and a Xilinx Field Programmable Gate Array (FPGA) all executing a set of 8 representative benchmarks that includes HPC applications, Convolutional Neural Networks (CNNs) for objects detection, and heterogeneous codes. We show that all the considered devices are vulnerable to thermal neutrons. For some devices, the probability for thermal neutrons to generate an error appears to be higher than the probability due to high energy neutrons.

¹FIT (Failures in Time) rate is a measure of the number of device failures in one billion (10^9) device-hours of operation.

The main contributions of this paper are: (1) an experimental evaluation of the probability for a high energy vs. thermal neutron to generate an error in modern computing devices; (2) an estimation of the thermal neutrons flux modification due to materials heavily present in a supercomputer room, based on homemade thermal neutrons detectors; (3) the evaluation, based on (1) and (2), of the contribution of thermal neutrons to the error rate of computing devices.

The remainder of the paper is organized as follows. Section 2 serves as a background and reviews previous work. Section 3 describes our evaluation methodologies. Section 4 presents the effects of thermal neutrons on DDR memories. Section 5 quantifies our experimental results, Section 6 presents the estimated FIT rates, and Section 7 concludes the paper.

2. BACKGROUND AND RELATED WORKS

This section serves as a background and related works on high energy and thermal neutrons effects on computing devices.

2.1 Motivation

Radiation is a known cause of upsets in computers [8]. The interaction of particles, primarily neutrons for terrestrial machines, with transistors can reverse the value of the bits stored in memory or create current spikes in logic operations. These faults can be masked with no effect on the system functionality, corrected by ECC (Error Correcting Code), create an undetected error known as Silent Data Corruption (SDC), or create a Detected Unrecoverable Error (DUE). The most serious of these effects are SDCs and DUEs. It is well known that thermal neutrons can affect electronic devices [5, 7]; however, only devices containing ^{10}B are considered susceptible to thermal neutrons. Approximately 20% of naturally occurring boron is ^{10}B with the rest primarily being ^{11}B . Depleted boron, where the ^{10}B content is low, is produced in the nuclear power industry but it is expensive in large quantities and generally not used in COTS parts. Previous generations of silicon chips used borophosphosilicate glass in the manufacturing process until it was shown to increase by $8\times$ the likelihood of upsets and replaced with glass not containing boron [9, 10]. Until recently the "boron problem" was considered a solved issue, however, as our experiments show, newer silicon chips seem to have re-introduced large amounts of boron back into the manufacturing process. Understanding how this change effects overall system reliability is the primary motivation for this work.

Recently, ^{10}B was found in the manufacturing process of COTS devices [11]. Some previous work has studied the sensitivity of SRAM and FPGA devices to thermal neutrons [12, 13, 14]. Weulersse et al. [15] compared the error rates of some memories (SRAM, CLB, caches) induced by thermal neutrons, 60MeV protons, and 14MeV neutrons. This preliminary study shows that the sensitivity to thermal neutrons ranges from $1.4\times$ to $0.03\times$ the high energy neutron one. While very interesting, these experiments were conducted on mem-

ory devices not typically used in HPC systems. In addition, many memory errors can be masked or detected through ECC and parity on HPC systems. Unfortunately, Weulersse et al. do not share details about the kind of errors observed during their experiments (single vs. multiple bit flips), preventing extrapolation of their results to HPC devices with ECC enabled.

Our work advances the knowledge on HPC reliability by considering the impact of thermal neutrons on the reliability of HPC devices. The radiation experiments were performed on devices executing representative applications under normal operational configurations (i.e., protection mechanisms enabled) to provide a realistic comparison between the error rates induced by high energy and thermal neutrons. Unlike previous publications, we perform both thermal and high energy neutrons experiments on exactly the same devices in the same conditions to limit comparison uncertainty. Furthermore, for the first time, we investigate through thermal neutron detector measurements, how modern data center construction and cooling systems designs influence the thermal neutron flux and the HPC system fault rates.

2.2 High Energy and Thermal Neutrons

High energy neutrons, or fast neutrons, are produced by the interaction of galactic cosmic rays with the atmosphere. Neutrons with energies that range from 1 to over 1,000 MeV are known to disturb the function of electronic devices and are considered a main cause of faults in terrestrial electronic devices [5, 4]. High energy neutrons primarily interact with silicon chips via elastic scattering which can deposit thousands of electron-Volts (eV, a standard unit of energy used in nuclear physics) of energy into a recoil nuclei. Neutron scattering may also produced secondary particles such as protons or alphas. All of these processes can free bound electrons in large enough quantities to alter the behavior of the circuits on a chip which may induce faulty behavior in one or more bits. Permanent damage can also occur due to the displacement of atoms within a chip. The flux of high energy neutrons in the atmosphere has been thoroughly studied since Hess' discovery [16, 17]. The flux is known to vary across the surface, as a consequence of the earth's magnetic field, and increases exponentially with altitude, reaching a maximum at about 60,000 ft. Under normal solar conditions, the fast neutron flux is almost constant for a given latitude, longitude, and altitude.

Thermal neutrons, or slow neutrons, are low energy neutrons (lower than 0.5 eV), produced by the moderation of high energy neutrons in materials or the emission of neutrons from nuclear decay. Incident high energy neutrons rain down as part of cosmic ray induced showers reaching thermal energies after 10-20 interactions. Thermal neutrons continue moving until they are either absorbed in a nuclear reaction, or decay (whilst stable in atomic nuclei, free neutrons have a half-life of about 10 minutes). When a thermal neutron is absorbed by ^{10}B , the resulting isotope decays, producing a Lithium

isotope and an alpha particle. Both the Lithium isotope and alpha particle can induce faults. The amount of boron in a particular computing device is proprietary information that is not disclosed by industry. The only definitive way to evaluate the thermal neutron sensitivity of a device is to expose it to thermal neutrons.

The flux of thermal neutrons, in contrast to high energy neutrons, can be difficult to predict as it strongly depends on the environmental conditions as well the presence of other materials (primarily hydrogen containing) in the device’s immediate surroundings (like concrete, water, a fuel tank, etc..) in addition to latitude, longitude, and altitude. Various authors have made calculations to evaluate thermal fluxes in realistic cases [7, 18, 19, 20]. As a result, when predicting the error rate caused by thermal neutrons, it is essential to measure rates in realistic settings.

We have built and deployed a neutron detector in order to have a precise understanding of the thermal neutron flux inside a representative data center. We measured the rates of thermal neutrons in the proximity of materials such as water, concrete, or plastic and demonstrate that cooling water, for instance, can increase the thermal neutron flux (and thus error rate) by up to 20%. In Section 6 we estimate the high energy vs. thermal neutrons error rate for two locations with known neutron fluxes and discuss the effects of environmental conditions (sunny and rainy day) and surrounding materials (concrete slab floors).

2.3 Supercomputer Cooling

One of the main challenges in designing HPC systems is the dissipation of heat. A modern supercomputer can push more than 750 watts per square foot which can easily overwhelm traditional cooling systems [21]. Today’s supercomputers consist of hundreds of computing racks (e.g., Summit uses 256 racks [22]), requiring specific room designs to optimize both cooling efficiency and ease of maintenance.

One notable and growing trend in data centers is the use of liquid cooling [23]. Eight of today’s Top10 supercomputers use some form of liquid cooling [24]. Liquid cooling is more efficient at heat removal than traditional air cooling and using it allows for an overall increase in performance and power efficiency. Traditional data centers may use 25% to 35% of their energy budget just for cooling. IBM chiller-less water cooling systems have been shown to reduce the cooling energy overhead to just 3.5% [25]. IBM has noted that using liquid cooling in can allow for a 34% increase in processor frequency which can increase system performance by approximately 33% [26, 27].

3. METHODOLOGY

To evaluate the contribution of thermal and high energy neutrons to the error rate of devices it is necessary to: (1) measure the probability that a neutron will generate a fault, and (2) estimate the flux of high energy and thermal neutrons where the device will operate. We measure (1) through accelerated neutron beams exper-

iments and estimate (2) using existing data as well as initial measurements of actual thermal neutron rates in an approximate setting.

In this section, we describe the devices and applications chosen to test the impact of high energy and thermal neutrons in modern computing devices reliability. We also detail the radiation experiments setup used for this work and describe the detector we used to measure the impact of materials in the thermal neutron flux.

3.1 Devices

We selected six devices for this study using different technologies and vendors to have an in-depth insight of thermal neutrons sensitivity on a breadth of modern devices. It is worth noting that both the fabrication process and the foundry can significantly impact the amount of ^{10}B in the device.

Intel Xeon Phi is an HPC accelerator that, even if recently announced as dismissed, powers some of the fastest supercomputers from the Top500 list [24]. The Xeon Phi tested is the coprocessor 3120A, which implements the *Knights Corner* architecture, and it is built using a 22nm **Intel’s 3-D Tri-gate technology**.

NVIDIA K20 is a GPU built with the *Kepler* architecture and is fabricated in a 28nm **TSMC standard CMOS technology**. This model is specially built for HPC systems and has 2496 CUDA cores divided across 15 Streaming Multiprocessors (SMs).

NVIDIA TitanX is a GPU built with the *Pascal* architecture and fabricated in a 16nm **TSMC FinFET**, it has 3584 CUDA cores split across 28 SMs.

NVIDIA TitanV is built with the *Volta* architecture and fabricated in a 12nm **TSMC FinFET**, it features 5120 CUDA cores divided into 80 SMs.

AMD Accelerated Processing Unit (APU) is a heterogeneous device that integrates CPU and GPU in the same chip sharing the same memory. The APU considered is the AMD A10 7890K Kaveri fabricated in a 28nm **SHP Bulk Process at Global Foundries**. This device includes 4 Steamroller CPU cores and a GCN architecture AMD Radeon R7 Series GPU containing 512 cores with 866MHZ each. We consider three APU configurations: CPU, GPU, and CPU+GPU.

FPGA is the Zynq-7000 designed by Xilinx using a 28nm **TSMC technology**. The FPGA is composed mainly of Configurable Logic Blocks (CLBs), Digital Signal Processor (DSP) blocks, and embedded memory blocks (BRAM).

3.2 Codes

The set of devices we consider covers a wide range of architectural and computational characteristics. Using the same code for each device would bias the reliability evaluation, in favor of the devices that are more efficient in executing the chosen code. To have a fair evaluation, then, we choose for each class of devices the codes that better fit with its computational characteristics. For Xeon Phi and GPUs we chose four codes representative of **HPC**: MxM, LUD, LavaMD, and HotSpot. We selected three **heterogeneous** codes specially made to

fully utilize the APU architecture: SC, CED, and BFS. Finally, on GPUs and FPGA we tested two **neural networks** to represent codes that have a significant impact on self-driven vehicles: YOLO and MNIST.

Matrix Multiplication (MxM) is representative of highly arithmetic compute-bound codes used in HPC and for features extraction in CNNs [24].

LUD is a linear algebra method that calculates solutions for a square system of linear equations, representative of highly compute-bound codes [28].

LavaMD simulates particle interactions using Finite Difference Methods [28]. LavaMD is compute bound, being mostly composed of dot-products.

Hotspot is representative of stencil solvers [28], it estimates the processor temperature using an architectural floor plan and simulated power measurements.

Stream Compaction (SC) is a memory-bound code used in databases and image processing applications. SC is composed of a data manipulation primitive that removes elements from an array.

Canny Edge Detection (CED) extracts information from images and reduce the amount of data to be processed. CPU and GPU concurrently work on different frames. The input frames are a subset of the Urban Dataset used for neural networks training [29].

Breadth First Search (BFS) is a search in graphs algorithms that performs non-uniform memory access widely used in GPS Navigation Systems. The input graph we select for our evaluation represents the highways of the Great Lakes area in the US [30].

YOLO is a Convolution Neural Network (CNN) used for object classification and detection [31].

Modified National Institute of Standards and Technology (MNIST) is a CNN used for classifying handwritten digits [32]. We have tested MNIST only on FPGAs as it is a minimal network that would not exercise sufficient resources on GPUs or Xeon Phis.

3.3 Radiation Experiments Setup

To evaluate the sensitivity of our devices to high energy and thermal neutrons we exposed the devices on two different beamlines at the ISIS spallation neutron source in the UK: ChipIR for high energy neutrons and ROTAX for thermal neutrons.

ChipIR [33] is the reference beamline dedicated to the irradiation of microelectronics and it features a high energy neutron spectrum, as similar as possible to the atmospheric one. The flux with neutron energy above 10 MeV is $5.4 \times 10^6 n/cm^2/s$, while the thermal component ($E < 0.5eV$) is $4 \times 10^5 n/cm^2/s$ [34].

ROTAX [35] is a general purpose beamline with a thermal neutron spectrum generating a flux of $2.72 \times 10^6 n/cm^2/s$. Here the thermalization is achieved by moderation of the neutrons using liquid methane.

The spectra of the two beamlines are compared in Figure 1 on a log-log scale where the fluxes are proportional to the areas under the curves. As Figure 1 suggests, most neutrons in ROTAX are thermals and most neutron in ChipIR are high energy one.

To evaluate the sensitivity to thermal and high energy

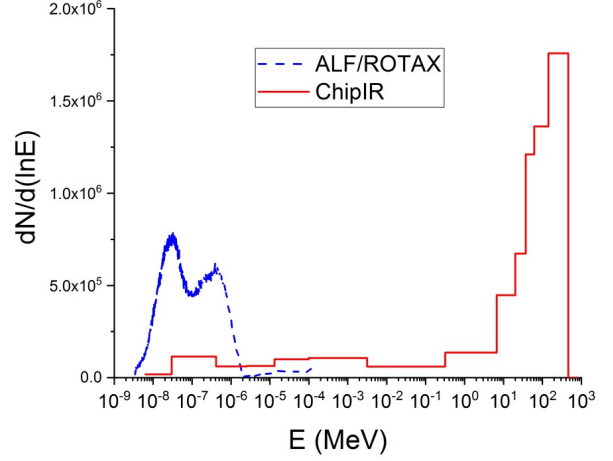


Figure 1: The neutron spectra of the beamlines used for irradiation in lethargy scale.

neutrons we align the devices described in Section 3.1 with the beam, while executing the codes listed in Section 3.2. The device output is compared with a pre-computed fault-free copy and any mismatch is marked as an SDC. If the application dies, gets stuck, or the device stops responding we count this as a DUE. Dividing the number of observed errors with the fluence the device has received we can calculate the device sensitivity, expressed as **cross section** [cm^2]. The higher the cross section, the higher the probability for one neutron (either thermal or high energy) to generate an observable error (either SDC or DUE).

To eliminate any setup-dependent differences between thermal and high energy neutrons, we irradiate the same physical devices executing the codes with the same input vector both in ROTAX and in ChipIR. It is worth noting that, apart from DDR that experienced permanent faults, testing the same device at ROTAX and then at ChipIR (or the other way around) does not influence the measured error rates. The only difference between the two experiments is that, thanks to the higher neutron energies, at ChipIR we can align various boards with the beam, as shown in Figure 2. Using a derating factor that takes distance into account we can measure the sensitivity of multiple devices in parallel. In ROTAX, as the irradiate devices stop most of the incoming thermal neutrons, we must test one device at a time. In Figure 3 we show the setup for the Titan V evaluation. Due to limitations in the thermal neutrons experiment, we could only test one sample of each device. The high energy neutrons error rate variation among different samples of the same device has already been shown to be low, recent works indicate a variation of about 10% [36, 37].

3.4 Thermal Neutrons Detector

We have designed and deployed a thermal neutron detector, called Tin-II, to measure the flux of thermal



Figure 2: Experimental setup in ChipIR. The arrow indicates the direction of the neutron beam.

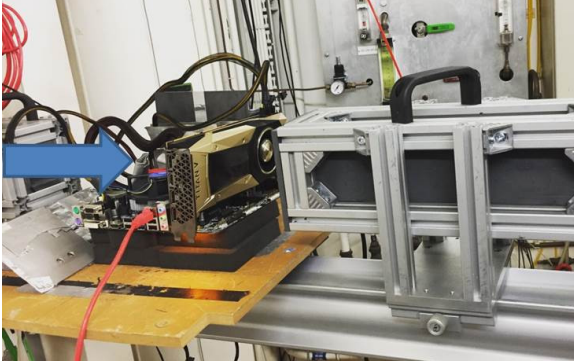


Figure 3: Titan X experimental setup in ROTAX. The arrow indicates the direction of the thermal neutron beam.

neutrons in different conditions. Ultimately, Tin-II will be used to measure the flux of thermal neutrons inside the data center housing the Trinity supercomputer at LANL. Tin-II consists of two identical ^3He cylindrical detectors. The interaction of radiation (neutrons, gammas, betas, etc.) with the detectors triggers a reaction that is amplified, filtered, and counted as an event.

We calibrated the two detectors for a period of 18 hours to ensure that they have the same detection efficiency. Then, we shielded one of the two cylinders with cadmium. Cadmium effectively blocks thermal neutrons, while being transparent to other types of radiation such as high energy neutrons, gammas, betas, etc. As a result, one of the two cylinder (bare detector) detects all radiation reactions, while the other (shielded detector) counts only radiation reactions that are not thermal neutrons. The difference in count rates between these two detectors, multiplied by an efficiency value, reflects the average thermal neutron flux.

Tin-II counted thermal neutron events over the course of several days. To estimate the effect of some of the characteristic materials in modern data centers on the thermal neutron flux, we placed a box containing 2 inches of water close to the detector. The count difference with and without the water, shown in details in Section 6.1, indicates its influence in the thermal neutrons flux.

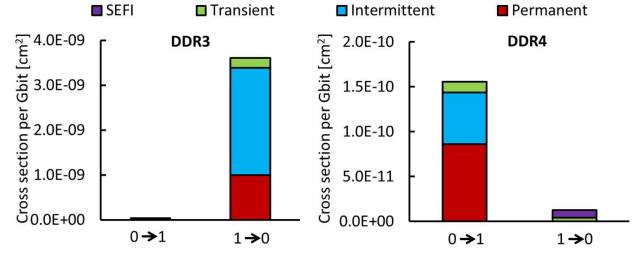


Figure 4: DDR3 and DDR4 thermal neutrons cross sections.

4. MEMORIES

In this section, we present the Double Data Rate (DDR3 and DDR4) Dynamic Random Access Memories (DRAM) sensitivity to thermal neutrons. Both DDR memories are synchronous DRAM tested without ECC and composed of a single rank x8 memory module. The DDR3 is a 4GB module that operates at 1.5V with a frequency of 1866 MHz and timings 10-11-10. The DDR4 is an 8GB module that operates at 1.2V with a frequency of 2133MHz and timings 13-15-15-28. As vendors are not explicitly mentioned, cross sections are shown in nominal values.

We irradiate the devices while performing a continuous read/write *correct loop*: banks are set to 0xFF (or 0x00) and continually read while irradiated with neutrons. When an unexpected value appears, error counters are increased, the corrupted data is downloaded for further analysis, and the memory bank is rewritten. This read/write loop allows differentiating 1-0 and 0-1 bit flips. While Static RAM has a symmetric structure, DDR are likely to be more sensitive to either one of the two possible bit flip directions (one-to-zero and zero-to-one), depending on the cell implementation and on the use of complementary logic.

The errors are classified into four categories:

- **Transient error:** a bit flip that does not systematically appear in the following memory read.
- **Intermittent error:** a memory location returns incorrect values, but not necessarily in consecutive reads. Intermittent errors have been seen in DDR and are dependent on environmental conditions, like temperature [38].
- **Permanent error:** a memory location consistently returns an incorrect value (stuck-at). Permanent errors are caused by Displacement Damage (the neutron dislocates atoms in the transistor) and can possibly be repaired with annealing (i.e., heating the device) [39, 40].
- **Single Event Functional Interrupt (SEFI):** a large portion of the memory array return incorrect values, likely caused by an error in the DDR control logic circuits. Further reads/writes will return correct values [41].

Figure 4 shows the thermal neutrons cross section per Gbit for DDR3 and DDR4. We do not report high energy neutron data since after few minutes of irradiation

at ChipIR both DDR3 and DDR4 experienced a high number of permanent faults, impeding further data collection. However, the sensitivity of DDR memories to high energy neutrons has been extensively studied, and experimental data can be found in [42, 43, 39, 44].

Figure 4 highlights that the DDR4 memory cross section is approximately one order of magnitude lower than the DDR3 one, showing significant reliability improvements probably resulting from new manufacturing processes as well as transistor placement enhancement. We also observe in Figure 4 that more than 95% of all the errors are in one of the two possible bit flip direction, one-to-zero for DDR3 and zero-to-one for DDR4. The opposite direction for DDR3 and DDR4 suggests that one device is manufactured with complementary logic. Another interesting point our data highlights is the proportion of each error category changes from DDR3 to DDR4. Permanent errors are more than 50% of all observed errors in DDR4, while on DDR3 only less than 30% of errors are permanent. It is also worth noting that both technologies present SEFI errors during the experiments. That is, an impinging particle on both DDR memory control circuits tend to incite similar malfunctioning behaviors.

Finally, all the observed transient and intermittent errors were single bit flip. This is a promising result, as SECDED ECC is shown to be sufficient to corrects most thermal neutrons induced errors [45]. On the contrary, in a SEFI error multiple corrupted bits were observed.

5. CROSS SECTION RESULTS

In this section, we compare the cross section measured at ChipIR and ROTAX for the tested devices and codes with the methodology described in Section 3.3. We emphasize that we used exactly the same device and setup for both ChipIR and ROTAX experiments. Due to beam time limitations (mainly at ROTAX as we must test only one device at a time) we could not test all the benchmarks in each device. Recall that a higher cross section indicates a higher probability of a single (high energy or thermal) neutron inducing faults. To evaluate the impact of thermal vs. high energy neutrons on the device error rate we need to consider the natural background flux, which is done in Section 6.

As we show, the cross section to thermal neutrons is far from being negligible, indicating the presence of ^{10}B in the silicon doping. Reported data have been normalized to the lowest cross section for each vendor to prevent the leakage of business-sensitive data while allowing a direct comparison between codes and devices of the same vendor. We also report error bars considering Poisson’s 95% confidence interval.

Figure 5 shows the **Xeon Phi** SDC and DUE cross sections for high energy and thermal neutrons. On average the thermal neutrons cross section is much lower (1/20) than the high energy neutrons’ one, for both SDC and DUE. This low sensitivity to thermal neutrons is a sign that either little boron is used in the production of Xeon Phi or depleted boron is used. HotSpot is the most sensitive code for both SDCs and DUEs.

HotSpot is especially sensitive to DUEs, with a cross section more than $2\times$ higher than the average for both high energy and thermal neutrons. HotSpot, in contrast to the other codes, uses a high number of control flow statements and has low arithmetic intensity, increasing the sensitivity to DUEs.

For SDCs, the high energy neutron cross sections vary significantly depending on the code being executed (more than $2\times$ across codes), which is in accordance with previous work [36, 46]. The SDC cross sections for thermal neutrons, however, have a very low variation between codes (less than 20%) which may be an artifact of the low number of SDCs observed. This result suggests there is a negligible sensitivity to thermals in the chip resources that are responsible for the variation between error rates in the high energy SDC results. DUEs, on the other hand, have a similar trend for high energy and thermal neutrons. DUE faults induced by thermal neutrons seem to have similar effects to DUE faults induced by high energy neutrons.

Figure 5 shows the sensitivity of **NVIDIA GPUs** to thermal and high energy neutrons. For the K20, on the average, both the SDCs and DUEs thermal cross sections are very high, being 60% and 50% of the high energy neutrons ones. This indicates the presence of a significant amount of ^{10}B in the manufacturing process. The thermal neutrons SDC cross section trend across codes is also similar to the high energy neutrons one, in the sense that the code with the largest thermal neutrons cross section (i.e., HotSpot) is also the code with the largest high energy neutron cross section. This suggests that ^{10}B is present in the computing resources and memory of these devices, and that the fault locations are similar for both kind of neutrons.

It is also interesting to notice that YOLOv2 is the only code for which DUEs are more likely than SDCs, for both kind of neutrons. This result follows previous work that shows low SDC sensitivity in CNN based object detection [47]. As shown in Figure 5, YOLOv2 DUE cross section for thermal neutrons is more than 50% higher than the DUE cross section for high energy neutrons and more than $2\times$ higher than the average of all K20 codes. This cross section indicates that the reliability for YOLOv2 in environments where thermal neutron flux is significant will be much worse than expected, especially for a safety-critical application like self-driven cars.

For Titan X and Titan V, on the average, the thermal neutron cross section is an order of magnitude lower than the high energy one. The impact of thermal neutrons is lower for the newest GPUs than on the mature K20. This may imply that FinFET based GPUs are less susceptible to thermal neutrons than CMOS GPUs (K20 is built using CMOS planar transistors, Titan X and Titan V using FinFET). However, for the MxM tests, Titan V (12nm) shows an almost doubled thermal neutron SDC cross section compared to the Titan X (16nm). Unfortunately, we were not able to test more codes on the Titan V and, at this point, we cannot confirm if the increased thermal neutron cross section is

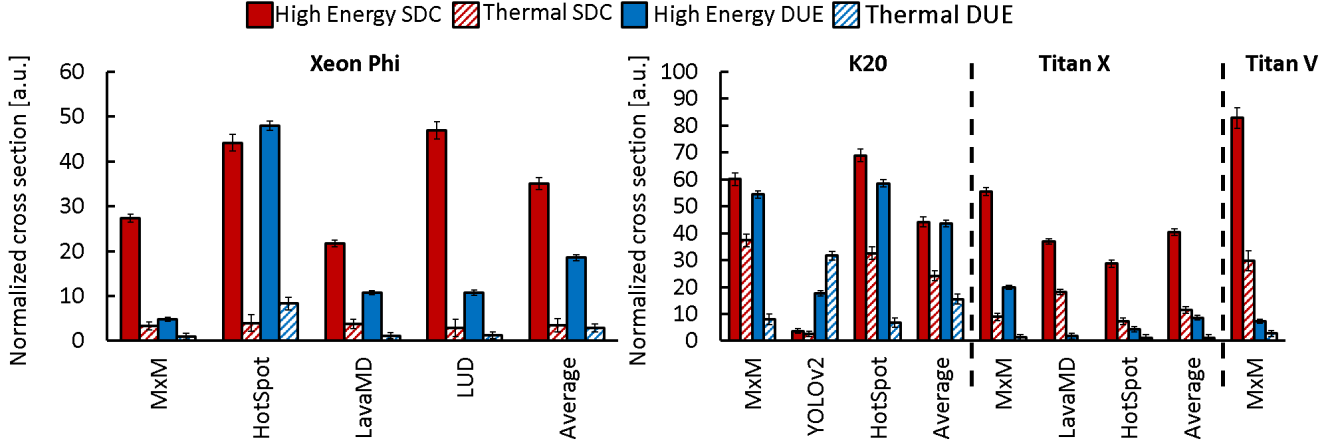


Figure 5: High energy and thermal neutrons normalized cross sections for Xeon Phi and GPUs.

intrinsic of smaller FinFET technologies.

The **AMD APU** cross sections are shown in Figure 6. As described in Section 3.1, the APU embeds a GPU and a CPU. We test the three heterogeneous codes described in Section 3.2 (CED, SC, and BFS) as executed on the GPU only, on the CPU only, and distributing concurrently 50% of the workload to the CPU and 50% to the GPU (CPU+GPU).

The APU-GPU, APU-CPU, and CPU+GPU SDC cross section for both thermals and high energy neutrons vary of more than an order of magnitude, forcing the use of logarithmic scale for APU data in Figure 6. The reported data shows that, on the average, the thermal neutrons cross section is reduced by between 1/4 and 1/5 the high energy neutron's, for CPU, GPU, and CPU+GPU. All APU configurations, on average, are more sensitive to SDCs than DUEs. It is also worth noting that the APU-CPU has, on average, a higher SDC sensitivity than APU-GPU. This is in accordance with previous work that shows a much lower probability for a fault in the AMD GPU to impact the application output than a fault in the CPU [48]. Additionally, in the APU, the GPU has a significantly smaller physical area than the CPU which reduces the probability of a neutron striking it and causing corruption.

A specific result to highlight is that SC code, which is the only memory-bound code of the three we test on the APU, has an SDC sensitivity to both high energy and thermal neutron extremely high when compared to others. As already shown, when the device is in idle waiting for data to be fetch from memory, registers and caches are exposed to radiation and store critical data [46]. Moreover, as observed for the Xeon Phi, the DUE cross section variation across different codes is much smaller than SDC variation. Finally, BFS has a particularly high DUEs sensitivity when the GPU is involved in computation (APU-GPU and CPU+GPU). This DUE increase is likely caused by the much higher stress in the CPU-GPU synchronization that BFS imposes by launching several GPU kernels (refer to section 3.2).

Figure 6 shows **Xilinx FPGA** SDC cross section

when executing the MNIST CNN. It is worth noting that neutron-induced errors in the configuration memory of SRAM FPGAs have a *persistent* effect, in the sense that a corruption changes the implemented circuit until a new bitstream is loaded in the device. The observation of an error at the FPGA output indicates that the bitstream has probably been corrupted. We reprogram the FPGA at each observed output error to avoid the collection of a stream of corrupted data, making the observation of DUEs very rare. In fact, as FPGA executes operation without any operating system, interfaces, or control-flow involved, a considerable amount of errors would need to accumulate in the configuration memory to have the circuit functionality compromised. We never observed a DUE in FPGAs during our experimental campaign.

We have tested two different versions of the neural network, one using double and the other using single precision floating-point arithmetic. When comparing the high energy and thermal neutrons cross sections for the two configurations, we can clearly perceive that the Xilinx FPGA is more sensitive to high energy neutrons. However, the thermal neutrons cross section is far from being negligible.

The double precision version takes about twice as many resources to be implemented in the FPGA. As the neutrons cross section is directly related to the circuit's area, the cross section is expected to be higher for the double version of MNIST. Experimental results for both high energy and thermal neutrons confirm this intuition. The thermal neutrons cross section for the double version is particularly higher than the single one, being almost four times larger.

Our results show that different codes executed on the same device can have very different high energy vs thermal neutrons sensitivities. The physical interaction of a thermal neutron and, consequently, the resulting fault model (i.e., the way the physical fault is manifested at circuit level) and the impact on the code execution is highly different from the high energy neutron one. High energy neutrons can interact with any atom in the chip

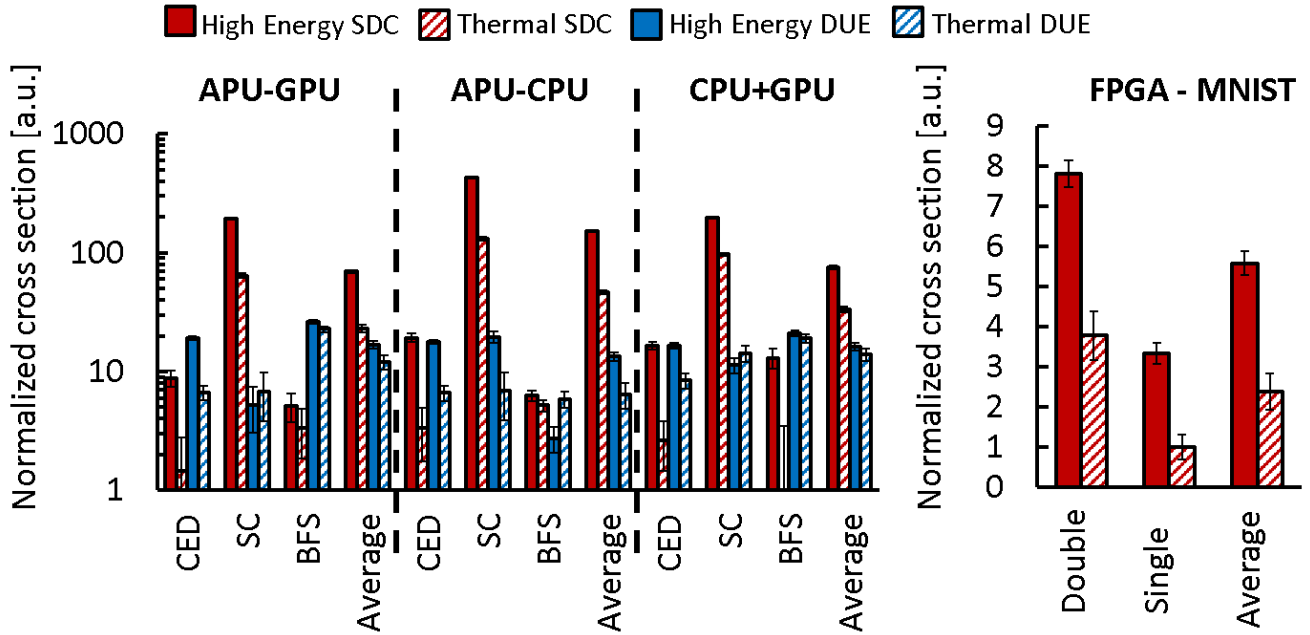


Figure 6: High energy and thermal neutrons normalized cross sections for AMD APU and FPGA.

or package material, triggering a reaction that may potentially reach a transistor’s vulnerable area. The fault can happen some distance from the high energy neutron impact, and the particles resulting from the interaction can travel in different directions [5]. Thermal neutrons, on the other hand, interacting only with ^{10}B , produce an alpha particle and lithium recoil that have very short ranges. When a thermal neutron is absorbed in virtually all other materials used in semiconductor manufacturing, the resulting nuclei typically only produces gamma rays, which do not produce bit flips.

Software fault-injection can emulate predefined fault models and study their effects, but cannot be used to study the fault manifestation nor to define different fault models. One way to investigate the different fault models would be to simulate the physical implementation of a transistor in a given technology and observe the effect of neutron strikes at different energies [49]. However, transistor implementation details are not available for COTS devices, which makes the comparison of the beam experiment cross sections of various codes the only possible way to highlight code-dependent thermal vs high energy neutrons induced error rates.

6. FIT RATE ANALYSIS

The cross sections reported and discussed in Section 5, represent the device’s sensitivity to thermal or high energy neutrons. To have an understanding of the impact of thermal and high energy neutrons in the device error rate, we need to consider also the natural background radiation fluxes of each. FIT rates can then be calculated by multiplying the experimentally measured cross sections by the neutron fluxes. For DDR, we show absolute FIT rates while for computing devices, to avoid the

leakage of business sensitive data, we only show in percentages the contribution of thermal and high energy neutrons to the device’s FIT rates. This information allows us to evaluate how much thermal neutrons increases the FIT of each device. This also tells us how much the FIT rate of each device is underestimated if thermal neutrons are not considered.

6.1 Thermal Neutrons Flux

The flux for high energy (fast) neutrons in the atmosphere can be precisely estimated considering the altitude, longitude, latitude, and solar activity using on-line available tools [50]. However, the environment and the materials that surround a device significantly impact neutron flux and energy. Materials such as concrete and water scatter neutrons which lose energy with each interaction. For instance, during thunderstorms the rain droplets act as moderators slowing high energy neutrons into lower energy ones. The thermal neutron flux, as measured in [7], can be as much as $2\times$ higher during a rain storm than on a sunny day. Thermal neutron rates may be as much as 20% higher over a large slab of concrete such as in a parking lot or the concrete floor of a machine room. Water cooling systems can also have the side effect of significantly increasing the proportion of thermal neutrons that strike a device.

In order to empirically measure the impact of materials in the thermal neutron flux in a data center, we placed the Tin-II detector (details in Section 3.4) in a building similar to the one containing the Trinity supercomputer. We collected data over the course of several days, then placed 2 inches of water in a pan over the detector starting on 20th April 2019. Figure 7. When water is placed over the detector the thermal neutron counts abruptly increase of about 24%. This increase

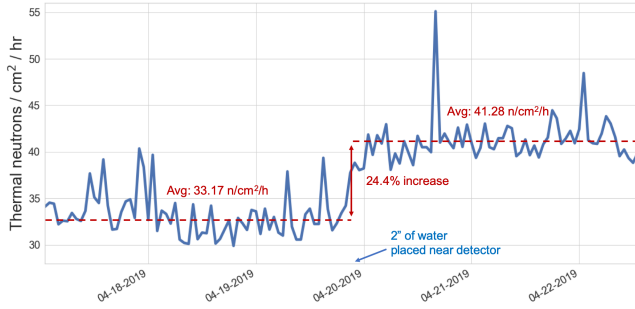


Figure 7: Tin-II thermal neutron detector measurements with two inches of water placed over detector on 20th April 2019.

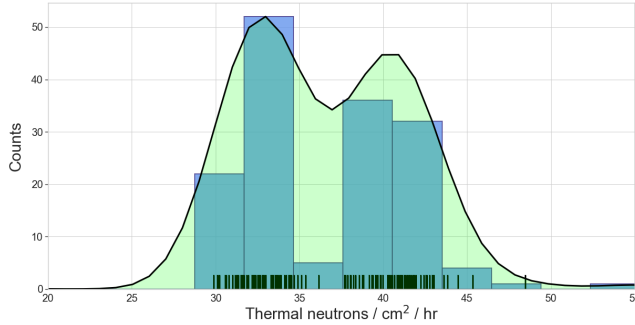


Figure 8: Tin-II thermal neutrons count rate. The bimodal distribution shows the influence of water placed above the detector.

shows that the presence of water in the cooling system can significantly increase the rates of thermal neutrons in a system, which in turn will increase the rates in the devices sensitive to those neutrons as seen in section 5.

Furthermore, to confirm the statistical significance of the thermal neutron flux calculation in our Tin-II experiments, Figure 8 displays the count rate of the thermal neutrons over time for our detector, the density function uses a Gaussian kernel density estimate. As can be seen, there are two main groupings of the data. The bimodal distribution is due to the water being placed over the detector on 20th April 2019.

The shape and placement of a water cooling system can impact the way thermal neutrons are produced. The LANL's Trinity supercomputer's water cooling pipes are below the machine which sits several feet above a concrete slab, whereas ORNL's Summit machine sits directly on a concrete slab with water cooling pipes running overhead. Both of these machines have liquid cooling radiators in the racks. Based on physical considerations, we believe the final flux for most liquid-cooled machines will be elevated. Figures 9 and 10 show simple simulations of the Trinity supercomputer for an incident neutron uniform distribution in the 1-15 MeV range modeled using the Los Alamos National Laboratory MCNP (Monte Carlo N-Particle) code coupled with ENDFVII neutron cross sections [51]. Figure 9 shows the distribution of thermal neutron flux in the

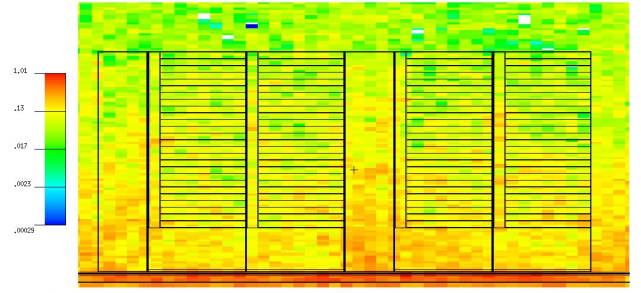


Figure 9: MCNP simulation of the distribution of the thermal neutron flux in racks of a Cray XC40 water cooled system (front view). Red indicates a higher rate in the lower blades of the rack as fast neutrons from above lose energy and "thermalize" while passing through the rack and cooling system.

racks of Trinity (front view) where each rack is composed of 26 computer blades. The simulated thermal neutron distribution per each blade is mapped in Figure 10.

These same considerations exist when trying to understand the thermal neutron component of faults in autonomous vehicles. The road material, concrete or asphalt, the vehicle is driving on makes a difference, as does the weather, and the type and volume of fuel the vehicle uses. In addition, the number of passengers will change the thermal neutron flux, as humans are primarily composed of water which makes us excellent neutron moderators.

6.2 High Energy vs Thermal Neutrons FIT

The average thermal neutron flux at New York City on a sunny day, excluding surrounding materials such as water or concrete, has been measured to be approximately $4n/cm^2/s$ [7]. Multiplying the flux by the DDR memory cross sections measured in Section 4 we can estimate the DDR3 and DDR4 thermal neutron-induced error rate in NYC to be about 3.09 and 0.14 FIT per GB, respectively. These rates can increase by $2\times$ or more because of the impact of environmental conditions on the thermal neutrons flux.

For computing devices, in Figure 11 we show the percentage of the total FIT rates due to high energy and thermal neutrons. These calculations use measured values of neutrons at sea level (NYC) and in Leadville, CO (10,151 ft in altitude). The thermal rates used have been adjusted to compensate for back scattered neutrons from a concrete slab and water cooling as measured by Tin-II detector, an overall increase of 44% in the thermal flux. Note that on a rainy day the thermal flux may be as much as doubled over the rates used in this graph and the corresponding FIT rate on those days will increase in a corresponding way [7].

Xeon Phi processors, as stated in Section 5, have a low sensitivity to thermals, which is a symptom of the use of either depleted boron or a reduction in boron usage. Thus, the thermals FIT rate seen in figure 11 is a

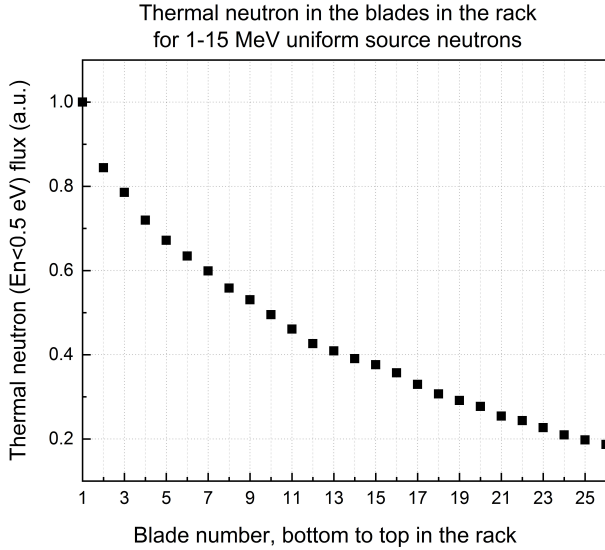


Figure 10: MCNP simulation of the distribution of the thermal neutron flux by "height" in a rack. Lower numbered blades are closer to the machine room floor and show a higher rate of thermal neutrons as fast neutrons from above slow down while passing through the rack.

relatively small percentage of the overall FIT rate (from 4.2% at NYC SDC up to 10.6% for Leadville DUE). The other tested devices, especially the K20 and CPU+GPU devices, have thermal FIT rates comparable to the FIT rates from high energy neutrons. At Leadville, K20 has 29% of the SDC FIT rate caused by thermal neutrons while APU CPU+GPU has 39% of DUEs caused by thermal neutrons.

6.3 Discussion

Figure 11 shows that if thermal neutrons contribution to the device error rate is not considered both the DUE and SDC FIT rates could be significantly underestimated, posing unconsidered risks to a safety critical application or reducing the HPC server productivity unexpectedly. Of particular interest in Figure 11 is the relatively high percentage of faults that result in Silent Data Corruption (SDC) on several of the tested devices. In general, HPC systems are designed and engineered to maintain SDC rates low and manageable, where corrupted calculations are rare and often noticeable to users. However, anything that increases the SDC rate is always concerning. In safety critical applications, SDCs should be strictly avoided as they could put the system in unexpected states, and they could potentially lead to unpredictable actions.

The elevated DUE rates are also of concern to HPC systems as they invariably result in a system crash and loss of some portion of a calculation's run time. It is worth noting that even with thin layers of shielding, embedded devices in vehicles can suffer from a much higher thermal flux than the one considered in Figure 11 due to moderation and reflection from the driver and

passengers, from cooling liquids, from ground and from the fuel tank filled with a hydrogen-rich fuel [52].

Our analysis shows that thermal neutrons are a threat for the reliability of supercomputers and safety critical applications that rely on COTS HPC devices. While the benefits in terms of cost, performances, and efficiency of COTS devices are not in question, their utilization in applications for which reliability is a concern must be coupled with a careful reliability evaluation that considers the impact of thermal neutrons. As the amount of ^{10}B in the manufacturing process is not publicly available, radiation experiments are one of the few ways to evaluate the sensitivity of a COTS device to thermal neutrons. Moreover, as the thermal neutron flux strongly depends on environmental conditions, the device error rate varies significantly when conditions change. Therefore it is critical to consider the realistic conditions in which the device will operate and estimate the correspondent thermal neutrons flux. These conditions have a direct impact on HPC applications. For instance, when supercomputer time is allocated, the checkpoint frequency may need to consider weather conditions. Dissimilarly to high-energy neutrons, thermal neutrons flux can be effectively reduced shielding the device with thin layers of cadmium or some inches of boron plastic. Unfortunately, cadmium is highly toxic and should not be heated, so it should not be placed in the proximity of an HPC device or of a cooling system, and boron plastic also thermally isolate the device, so it is impractical to be used as a shield between the cooling system (one of the most efficient sources of thermal neutrons) and the device.

7. CONCLUSIONS

In this paper we have experimentally investigate the differences between high energy and thermal neutron induced error rates in modern HPC devices. While purifying the Silicon dopant to remove ^{10}B would make devices immune to thermal neutrons, most COTS still use natural Boron. By irradiating devices with high energy and thermal neutrons while executing representative applications, we have demonstrated that thermals significantly impact device reliability. We have demonstrated that the impact of high energy and thermal neutrons depends not only on the specifics of the hardware, but also on the executed code. The impinging neutron energy has more or less effect depending on how the code accesses memory and executes instructions.

We have also shown that the FIT rates can vary based on the physical layout of the machine room in which a system resides and variations such as weather conditions external to the building.

The reported data attests the importance of thermal neutron reliability evaluation, which can significantly raise the total device error rate. As a future work, we plan to irradiate with thermal and high energy neutrons specific resources or components to deeply investigate different fault models. We also plan a thorough and sophisticated modeling of one or more data centers as well as the effects of different cooling regimes.

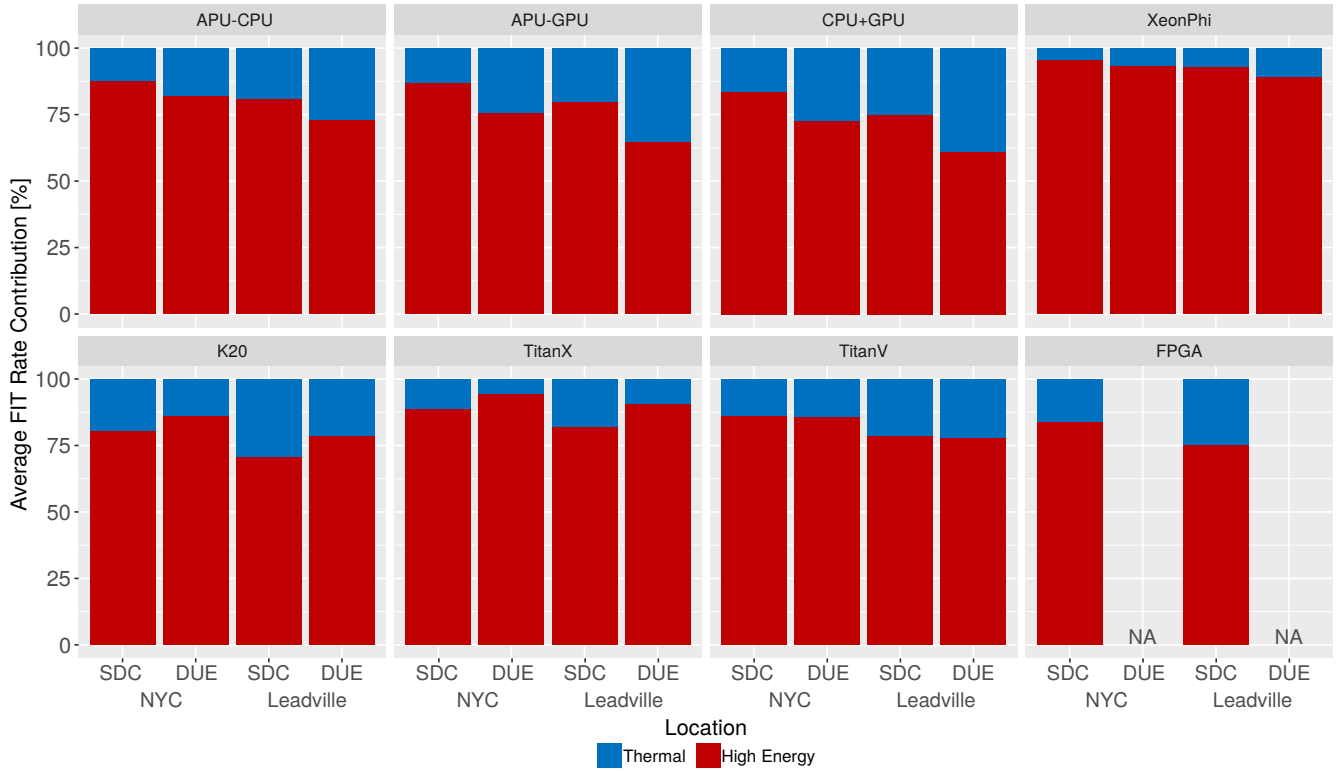


Figure 11: Percentage of total FIT rate due to high energy and thermal neutrons. All tested parts except Xeon Phi show significant errors due to ^{10}B levels.

8. REFERENCES

- [1] R. Lucas, "Top ten exascale research challenges," in *DOE ASCAC Subcommittee Report*, 2014.
- [2] A. Cohen, X. Shen, J. Torrellas, J. Tuck, Y. Zhou, S. Adve, I. Akturk, S. Bagchi, R. Balasubramonian, R. Barik, M. Beck, R. Bodik, A. Butt, L. Ceze, H. Chen, Y. Chen, T. Chilimbi, M. Christodorescu, J. Criswell, C. Ding, Y. Ding, S. Dwarkadas, E. Elmroth, P. Gibbons, X. Guo, R. Gupta, G. Heiser, H. Hoffman, J. Huang, H. Hunter, J. Kim, S. King, J. Larus, C. Liu, S. Lu, B. Lucia, S. Maleki, S. Mazumdar, I. Neamtiu, K. Pingali, P. Rech, M. Scott, Y. Solihin, D. Song, J. Szefer, D. Tsafir, B. Urgaonkar, M. Wolf, Y. Xie, J. Zhao, L. Zhong, and Y. Zhu, "Inter-disciplinary research challenges in computer systems for the 2020s," tech. rep., National Science Foundation, USA, 2018.
- [3] J. Dongarra, H. Meuer, and E. Strohmaier, "ISO26262 Standard," 2015.
- [4] J. Ziegler and H. Puchner, *SER-history, Trends and Challenges: A Guide for Designing with Memory ICs*. Cypress, 2004.
- [5] R. Baumann, "Radiation-induced soft errors in advanced semiconductor technologies," *Device and Materials Reliability, IEEE Transactions on*, vol. 5, pp. 305–316, Sept 2005.
- [6] M. Snir, R. W. Wisniewski, J. A. Abraham, S. V. Adve, S. Bagchi, P. Balaji, J. Belak, P. Bose, F. Cappello, B. Carlson, et al., "Addressing failures in exascale computing," *International Journal of High Performance Computing Applications*, p. 1094342014522573, 2014.
- [7] J. Dirk, M. E. Nelson, J. F. Ziegler, A. Thompson, and T. H. Zabel, "Terrestrial thermal neutrons," *IEEE Transactions on Nuclear Science*, vol. 50, no. 6, pp. 2060–2064, 2003.
- [8] JEDEC, "Measurement and Reporting of Alpha Particle and Terrestrial Cosmic Ray-Induced Soft Errors in Semiconductor Devices," Tech. Rep. JESD89A, JEDEC Standard, 2006.
- [9] R. Baumann, T. Hossain, E. Smith, S. Murata, and H. Kitagawa, "Boron as a primary source of radiation in high density drams," in *1995 Symposium on VLSI Technology. Digest of Technical Papers*, (Kyoto, Japan, Japan), pp. 81–82, IEEE, IEEE, 1995.
- [10] E. Normand, K. Vranish, A. Sheets, M. Stitt, and R. Kim, "Quantifying the double-sided neutron seu threat, from low energy (thermal) and high energy >10 mev) neutrons," *IEEE transactions on nuclear science*, vol. 53, no. 6, pp. 3587–3595, 2006.
- [11] S.-J. Wen, S. Pai, R. Wong, M. Romain, and N. Tam, "B10 finding and correlation to thermal neutron soft error rate sensitivity for srams in the sub-micron technology," in *2010 IEEE International Integrated Reliability Workshop Final Report*, pp. 31–33, IEEE, 2010.
- [12] S. Lee, I. Kim, S. Ha, C.-s. Yu, J. Noh, S. Pae, and J. Park, "Radiation-induced soft error rate analyses for 14 nm finfet sram devices," in *2015 IEEE International Reliability Physics Symposium*, pp. 4B–1, IEEE, IEEE, 2015.
- [13] Y.-P. Fang and A. S. Oates, "Characterization of single bit and multiple cell soft error events in planar and finfet srams," *IEEE Transactions on Device and Materials Reliability*, vol. 16, no. 2, pp. 132–137, 2016.
- [14] P. Maillard, M. Hart, J. Barton, P. Jain, and J. Karp, "Neutron, 64 mev proton, thermal neutron and alpha single-event upset characterization of xilinx 20nm ultrascale kintex fpga," in *2015 IEEE Radiation Effects Data Workshop (REDW)*, pp. 1–5, IEEE, 2015.
- [15] C. Weulersse, S. Houssany, N. Guibaud, J. Segura-Ruiz, J. Beaucour, F. Miller, and M. Mazurek, "Contribution of thermal neutrons to soft error rate," *IEEE Transactions on Nuclear Science*, vol. 65, no. 8, pp. 1851–1857, 2018.

- [16] V. F. Hess, "Über den ursprung der durchdringenden strahlung," *Z. Phys.*, vol. 14, p. 610, 1913.
- [17] J. F. Ziegler, "Terrestrial cosmic rays," *IBM Journal of Research and Development*, vol. 40, pp. 19–39, Jan 1996.
- [18] A. Hands, P. Morris, K. Ryden, C. Dyer, P. Truscott, A. Chugg, and S. Parker, "Single event effects in power mosfets due to atmospheric and thermal neutrons," *IEEE Transactions on Nuclear Science*, vol. 58, no. 6, pp. 2687–2694, 2011.
- [19] R. Baumann, "Soft error characterization and modeling methodologies at texas instruments," in *Proc. Semiconductor Research Council 4th Topical Conf. Reliability*. [CD-Rom] SemaTech CD-ROM, (USA), pp. 0043–3283, SemaTech, 2000.
- [20] R. Sheu and S. Jiang, "Cosmic-ray-induced neutron spectra and effective dose rates near air/ground and air/water interfaces in taiwan," *Health physics*, vol. 84, no. 1, pp. 92–99, 2003.
- [21] M. K. Patterson and D. Fenwick, "The state of datacenter cooling," *Intel Corporation White Paper. Available at <http://download.intel.com/technology/eep/data-center-efficiency/stateof-date-center-cooling.pdf>*, 2008.
- [22] "Summit system overview." https://www.olcf.ornl.gov/wp-content/uploads/2018/05/Intro_Summit_System_Overview.pdf, 2006.
- [23] A. Capozzoli and G. Primiceri, "Cooling systems in data centers: state of art and emerging technologies," *Energy Procedia*, vol. 83, pp. 484–493, 2015.
- [24] J. Dongarra, H. Meuer, and E. Strohmaier, "TOP500 Supercomputer Sites: November 2018," 2018.
- [25] T. Gao, M. David, J. Geer, R. Schmidt, and B. Sammakia, "Experimental and numerical dynamic investigation of an energy efficient liquid cooled chiller-less data center test facility," *Energy and buildings*, vol. 91, pp. 83–96, 2015.
- [26] M. Ellsworth, L. Campbell, R. Simons, M. Iyengar, R. Schmidt, and R. Chu, "The evolution of water cooling for ibm large server systems: Back to the future," in *2008 11th Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems*, pp. 266–274, IEEE, IEEE, 2008.
- [27] M. J. Ellsworth, G. F. Goth, R. J. Zoodsma, A. Arvelo, L. A. Campbell, and W. J. Anderl, "An overview of the ibm power 775 supercomputer water cooling system," *Journal of Electronic Packaging*, vol. 134, no. 2, p. 020906, 2012.
- [28] S. Che, M. Boyer, J. Meng, D. Tarjan, J. W. Sheaffer, S.-H. Lee, and K. Skadron, "Rodinia: A benchmark suite for heterogeneous computing," in *Proceedings of the IEEE International Symposium on Workload Characterization (IISWC)*, (Austin, TX, USA), pp. 44–54, IEEE, 2009.
- [29] K. Fragkiadaki, W. Zhang, G. Zhang, and J. Shi, "Two-granularity tracking: Mediating trajectory and detection graphs for tracking under occlusions," in *European Conference on Computer Vision*, pp. 552–565, Springer, Springer, 2012.
- [30] DIMACS, "9th dimacs." www.dis.uniroma1.it/challenge9/download.shtml, 2006.
- [31] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," *CoRR*, vol. abs/1506.02640, 2015.
- [32] L. Deng, "The mnist database of handwritten digit images for machine learning research [best of the web]," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [33] C. Cazzaniga and C. D. Frost, "Progress of the scientific commissioning of a fast neutron beamline for chip irradiation," in *Journal of Physics: Conference Series*, vol. 1021, p. 012037, IOP Publishing, IOP, 2018.
- [34] D. Chiesa, M. Nastasi, C. Cazzaniga, M. Rebai, L. Arcidiacono, E. Previtali, G. Gorini, and C. D. Frost, "Measurement of the neutron flux at spallation sources using multi-foil activation," *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, 2018.
- [35] H. Tietze, W. Schmidt, and R. Geick, "Rotax, a spectrometer for coherent neutron inelastic scattering at isis," *Physica B: Condensed Matter*, vol. 156, pp. 550–553, 1989.
- [36] D. Oliveira, L. Pilla, N. DeBardeleben, S. Blanchard, H. Quinn, I. Koren, P. Navaux, and P. Rech, "Experimental and analytical study of xeon phi reliability," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '17, (New York, NY, USA), pp. 28:1–28:12, ACM, 2017.
- [37] D. Oliveira, L. Pilla, M. Hanzich, V. Fratin, F. Fernandes, C. Lunardi, J. Cela, P. Navaux, L. Carro, and P. Rech, "Radiation-Induced Error Criticality in Modern HPC Parallel Accelerators," in *Proceedings of 21st IEEE Symp. on High Performance Computer Architecture (HPCA)*, ACM, 2017.
- [38] C. Constantinescu, "Intermittent faults and effects on reliability of integrated circuits," in *Reliability and Maintainability Symposium, 2008. RAMS 2008. Annual*, (Las Vegas, NV, USA), pp. 370–374, IEEE, IEEE, 2008.
- [39] H. Quinn, P. Graham, and T. Fairbanks, "Sees induced by high-energy protons and neutrons in sdram," in *2011 IEEE Radiation Effects Data Workshop*, pp. 1–5, July 2010.
- [40] J. R. Srour, C. J. Marshall, and P. W. Marshall, "Review of displacement damage effects in silicon devices," *IEEE Transactions on Nuclear Science*, vol. 50, pp. 653–670, June 2003.
- [41] E. I. Association *et al.*, "Test procedures for the measurement of single-event effects in semiconductor devices from heavy ion irradiation," *EIA/JEDEC Standard*, no. 57, 1996.
- [42] C. Constantinescu, "Impact of deep submicron technology on dependability of vlsi circuits," in *Dependable Systems and Networks, 2002. DSN 2002. Proceedings. International Conference on*, (Washington, DC, USA), pp. 205–209, IEEE, IEEE, 2002.
- [43] V. Sridharan, J. Stearley, N. DeBardeleben, S. Blanchard, and S. Gurumurthi, "Feng shui of supercomputer memory: positional effects in dram and sram faults," in *Proceedings of SC13: International Conference for High Performance Computing, Networking, Storage and Analysis*, p. 22, ACM, 2013.
- [44] S. M. Guertin and M. Cui, "See test results for the snapdragon 820," in *2017 IEEE Radiation Effects Data Workshop (REDW)*, pp. 1–6, July 2017.
- [45] V. Sridharan and D. Liberty, "A study of dram failures in the field," in *High Performance Computing, Networking, Storage and Analysis (SC), 2012 International Conference for*, pp. 1–11, IEEE, 2012.
- [46] V. Fratin, D. Oliveira, C. Lunardi, F. Santos, G. Rodrigues, and P. Rech, "Code-dependent and architecture-dependent reliability behaviors," in *2018 48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, pp. 13–26, IEEE, IEEE, 2018.
- [47] F. F. d. Santos, P. F. Pimenta, C. Lunardi, L. Draghetti, L. Carro, D. Kaeli, and P. Rech, "Analyzing and increasing the reliability of convolutional neural networks on gpus," *IEEE Transactions on Reliability*, vol. 68, pp. 663–677, June 2019.
- [48] H. Jeon, M. Wilkening, V. Sridharan, S. Gurumurthi, and G. H. Loh, "Architectural vulnerability modeling and analysis of integrated graphics processors," in *IEEE 10th Workshop on Silicon Errors in Logic - System Effects (SELSE)*, IEEE, 2013.
- [49] P. E. Dodd, "Physics-based simulation of single-event effects," *IEEE Transactions on Device and Materials Reliability*, vol. 5, pp. 343–357, Sep. 2005.
- [50] "Soft-error testing resource.." <http://www.seutest.com/cgi-bin/FluxCalculator.cgi>, 2006.

- [51] C. W. et al., ““mcnp6.2 release notes”,” 2018.
- [52] W. R. Leo, *Techniques for nuclear and particle physics experiments: a how-to approach*. Springer Science & Business Media, 2012.