

CRAFT: Criticality-Aware Fault-Tolerance Enhancement Techniques for Emerging Memories-Based Deep Neural Networks

Thai-Hoang Nguyen, *Student Member, IEEE*, Muhammad Imran, *Member, IEEE*, Jaehyuk Choi, *Member, IEEE*, and Joon-Sung Yang, *Senior Member, IEEE*

I. INTRODUCTION

Abstract—Deep Neural Networks (DNNs) have emerged as the most effective programming paradigm for computer vision and natural language processing applications. With the rapid development of DNNs, efficient hardware architectures for deploying DNN-based applications on edge devices have been extensively studied. Emerging Non-Volatile Memories (NVMs), with their better scalability, non-volatility and good read performance, are found to be promising candidates for deploying DNNs. However, despite the promise, emerging NVMs often suffer from reliability issues such as stuck-at faults, which decrease the chip yield/memory lifetime and severely impact the accuracy of DNNs. A stuck-at cell can be read but not reprogrammed, thus, stuck-at faults in NVMs may or may not result in errors depending on the data to be stored. By reducing the number of errors caused by stuck-at faults, the reliability of a DNN-based system can be enhanced. This paper proposes CRAFT, i.e., Criticality-Aware Fault-Tolerance Enhancement Techniques to enhance the reliability of NVM-based DNNs in the presence of stuck-at faults. A data block remapping technique is used to reduce the impact of stuck-at faults on DNNs accuracy. Additionally, by performing bit-level criticality analysis on various DNNs, the critical-bit positions in network parameters that can significantly impact the accuracy are identified. Based on this analysis, we propose an encoding method which effectively swaps the critical bit positions with that of non-critical bits when more errors (due to stuck-at faults) are present in the critical bits. Experiments of CRAFT architecture with various DNN models indicate that the robustness of a DNN against stuck-at faults can be enhanced by up to 10^5 times on CIFAR-10 dataset and up to 29 times on ImageNet dataset with only a minimal amount of storage overhead i.e., 1.17%. Being orthogonal, CRAFT can be integrated with existing fault-tolerance schemes to further enhance the robustness of DNNs against stuck-at faults in NVMs.

Index Terms—Deep learning hardware, Emerging Memories, Fault-Tolerance, Neural Networks, Stuck-at Faults

T. H. Nguyen is with the Department of Electrical and Computer Engineering, Sungkyunkwan University, Suwon 16419, Korea.

M. Imran is with the Department of Electrical Engineering, School of Electrical Engineering and Computer Science (SECS), National University of Sciences and Technology (NUST), Islamabad 44000, Pakistan.

J. Choi is with the Department of Semiconductor Systems Engineering, Sungkyunkwan University, Suwon 16419, Korea.

J.-S. Yang is with the School of Electrical and Electronic Engineering and Department of System Semiconductor Engineering, Yonsei University, Seoul 03722, Korea (e-mail: js.yang@yonsei.ac.kr)

This work was supported in part by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea Government (MSIT) under Grant 2022-0-00971; in part by the Next Generation Intelligent Semiconductor Development by the Ministry of Trade, Industry and Energy (MOTIE) under Grant 20011074; and in part by the Basic Science Research Program through the National Research Foundation of Korea funded by the Ministry of Education under Grant NRF-2020M3F3A2A01082326. The EDA Tools used in this work were supported by IDEC, Daejeon, South Korea.

DEEP Neural Networks (DNNs), a subset of Machine learning (ML) algorithms, have demonstrated impressive effectiveness in various applications such as computer vision, natural language processing, big data analysis and etc. A typical DNN consists of multiple hidden layers sandwiched between an input layer and an output layer. This hierarchical design allows DNN to solve complex programming tasks that appear to be infeasible with conventional programming approaches. However, despite the potential, DNNs often require enormous amount of computational power and hardware overhead, which makes it difficult to deploy them in real-time computing applications often running on mobile devices. As a result of the rapid development of DNNs, there is an enormous increase in the demand for efficient and scalable hardware architectures for DNNs' deployment. Realizing the high computational cost of DNNs, various methodologies have been proposed to achieve hardware-efficient architectures for DNNs [1], [2]. These techniques often focus on reducing the storage required by DNNs through network compression [1] and precision reduction [2]. Such methods have proven to be efficient, however, DNNs often need to sacrifice accuracy in exchange for a reduced implementation cost in resource-constrained devices.

Memory plays a key role in the applications involving large amount of data like DNNs. Current charge-based memory technologies such as Dynamic Random-Access Memory (DRAM), Static Random-Access Memory (SRAM) and Flash are facing challenges in continuing technology scaling [4]. Moreover, as the technology scales down, conventional memory technologies become highly prone to charge leakage which makes them a less attractive choice for data-intensive DNN applications. To cope with these issues posed by the conventional technologies, several emerging non-volatile memory technologies (NVMs) such as Resistive Random-Access Memory (ReRAM) and Phase Change Memory (PCM) have been extensively investigated over the past decade. With better scaling potential, better read performance and non-volatility [5], emerging NVMs are considered to be the potential replacement of the current charge-based memory technologies.

Beside being used for storage, thanks to their analog characteristic, emerging NVMs have also played a major role in designing high-performance and energy-efficient In-memory Computing (ICM) based accelerators for DNNs [6]–[8]. Such accelerators use emerging NVMs cells (i.e., ReRAM, PCM) to store the network's parameters and perform the matrix-vector

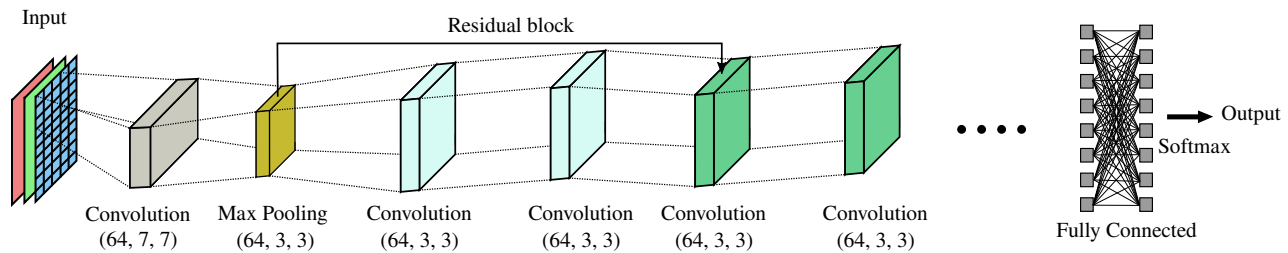


Fig. 1. Typical Convolutional Neural Network architecture (ResNet-18 [3])

multiplication in-place by organizing NVMs cells in a crossbar manner. With in-place computations, NVMs based IMC architecture eliminates the data movement between memory and separate computing units, which is found to be very costly in conventional von Neumann architectures. These features make the emerging memories an ideal choice for the future hardware implementations of DNNs.

Despite their promising features, emerging NVMs often suffer from hard error (i.e., stuck-at faults) [9]–[11] due to their low endurance and immature manufacturing process. A stuck-at fault occurs when the resistance/conductance state of an emerging NVMs cell can not be changed by a write operation. A stuck-at cell can still be read but not reprogrammed, thus, errors caused by stuck-at faults only arise when the stuck-at cell's state is not aligned with the desired data. Building on this insight, several fault-tolerance techniques have been proposed to increase the lifetime of an emerging NVMs-based memory system [12]–[16]. In NVMs-based DNN architectures, despite the inherent fault-tolerance of DNNs, a small number of stuck-at NVMs cells (especially those corresponding to the critical bits) can still cause a catastrophic loss to DNN's accuracy [17], [18]. Therefore, it is necessary to develop effective fault-tolerance enhancement techniques to mitigate such errors in NVMs-based DNN architectures.

Existing works on tolerating stuck-at faults in emerging NVMs have aimed for neuromorphic applications [19]–[23]. Despite being effective, these techniques often rely on an expensive retraining process of DNNs or a utilization of frequent auxiliary bits leading to a high storage overhead. On the other hand, several architectural techniques have also been proposed to tackle the problem of stuck-at errors in emerging NVMs, in general [12]–[16]. Such techniques also require a large amount of hardware storage and complex encoding/decoding mechanisms, making them infeasible for resource-constrained hardware with real-time performance requirements. To address the problems of previous existing works, we propose multiple lightweight yet effective techniques, collectively named CRAFT, to tolerate errors caused by the stuck-at faults in NVMs based DNN architectures. The first technique, called Intra-Block Address Remapping, effectively remaps the weights inside a block of data so that the impact of stuck-at faults on DNN's accuracy is minimized. The second method addresses the problem of single-bit error by simply inverting the data in the data block. Results of these two techniques have been presented in our earlier work [24]. To further enhance the robustness of the NVMs based DNNs, a novel Criticality-Aware Bits Switching method is proposed

which further enhances the DNN's accuracy in the presence of stuck-at faults by addressing the bit criticality in DNNs.

Rest of the paper is organized as follows. Section II covers the background of DNNs, emerging NVMs and stuck-at faults in emerging NVMs. Related works are presented in Section III. Section IV introduces the proposed Criticality-Aware Fault-Tolerance Enhancement Techniques (CRAFT). Finally, we evaluate the effectiveness of the CRAFT against existing techniques in Section V. Section VI concludes the paper.

II. BACKGROUND

A. Deep Neural Networks (DNNs)

Artificial neural networks (ANNs) are the computer algorithms inspired by the biological brain of animals. A layer of an ANN often consists of multiple nodes (i.e., neurons) connected to next layer through multiple connections (i.e., synapses/weights). Typical ANNs are made up of an input layer, an output layer and multiple hidden layers in between. A subset of ANN, Deep Neural Network (DNN), is an ANN with a large number of hidden layers (hence the name "Deep" Neural Network). Over the last decade, DNNs have made major breakthroughs in various fields, rendering the conventional programming approaches in these domains as obsolete. Especially, in the field of computer vision, Convolutional Neural Networks (CNNs) [3], [25] have attracted a lot of interest due to their exceptional effectiveness.

Fig. 1 shows a typical CNN (ResNet-18) architecture. A CNN is often composed of three types of layers : Convolutional, Fully Connected (FC), Pooling and Normalization. The convolutional layer is often used for extracting features of the input data by convolving the input with multiple relatively small size filters. The output of the convolutional layer is then fed into the pooling layer to reduce the spatial size of the representation. In ResNet architecture, as shown in the figure, the output data is propagated through multiple residual blocks consisting of two convolutional layers and a shortcut connection. Such blocks allow the CNN to increase its depth (i.e., numbers of layer) while preventing the vanishing/exploding gradient effect. At the end of the network, data undergoes a fully connected layer for classification followed by a softmax layer which outputs the probability of each class.

The hierarchical structure of DNNs allows them to outperform the conventional programming algorithms in solving complex problems by breaking them into simpler ones. However, DNNs require immense hardware resources for storing parameters and performing computations. This makes it extremely challenging to deploy a large-scale DNN on

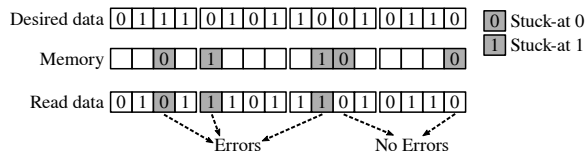


Fig. 2. Example of Stuck-at faults in emerging NVMs

resource-constrained hardware like mobile devices. To tackle this challenge, several techniques have been proposed to reduce the network size, thus making DNNs easier to deploy [1], [26]. The cost of using these techniques is the accuracy loss of DNNs. Depending on application, the accuracy loss may or may not be acceptable. In parallel with these approaches, several researches have investigated emerging non-volatile memories (NVMs) to provide high bandwidth, storage density and non-volatility for DNNs deployment. With such advantages over traditional charge-based memories, emerging NVMs are seen as ideal candidates for efficient and high-performance DNNs applications.

B. Emerging Non-Volatile Memories (NVMs)

Prominent emerging NVMs that are well-suited for DNN-based applications include Phase Change Memory (PCM) and Resistive RAM (ReRAM) [10]. Phase Change Memory consists of a chalcogenide phase-change material (e.g., $\text{Ge}_2\text{Sb}_2\text{Te}_3$) sandwiched between two electrodes. Data can be stored in PCM by modulating the phase-change material's state which is either crystalline or amorphous. The PCM cell has a low resistance in the crystalline state and a high resistance in amorphous state. ReRAM, on the other hand, consists of a conducting material (typically HfO_2) placed in between two electrodes [27]. The resistance state of a ReRAM cell can be changed by altering the concentration of defects in the conductive filament. A high defects concentration toward the bottom electrode changes the state of the ReRAM cell to the low-resistance state and a high defects concentration towards the top electrode leads to the high-resistance state. Both PCM and ReRAM have promising features of non-volatility, high switching speed and better endurance compared to existing Flash memory [10]. However, due to certain intrinsic characteristics of the underlying technology and an immature manufacturing process, these memories often face reliability issues such as hard errors (stuck-at faults) [12], resistance drift [?], [28]–[31] and write disturbance [19], [32]. This poses a challenge when employing emerging NVMs for DNNs applications.

C. Stuck-at Faults and DNNs Accuracy

1) *Stuck-at Faults in Emerging NVMs*: Stuck-at faults are a type of hard faults in emerging NVMs where the resistance state of a NVMs cell is locked at a certain state. When the cell resistance is fixed at the low resistance state, the fault is regarded as Stuck-at-Zero (SA0), otherwise, if the cell resistance is stuck at the high resistance state, this fault is considered to be Stuck-at-One (SA1). Depending on the data to be stored and the stuck-at state, a stuck-at fault may or may

not cause error in the system. Fig. 2 depicts a phenomenon of stuck-at faults in emerging NVMs. The first row shows the correct data expected to be stored and read from the memory. Second row shows the location and state of the stuck-at faults in emerging NVMs and the third row indicates the erroneous data read from the memory. As shown in the figure, the last two stuck-at fault locations do not introduce any error in the data because these cells' desired data is in-line with their stuck-at states. On the other hand, a mismatch between desired data and the stuck-at state causes error, which unintentionally flips the corresponding bit. Therefore, by aligning the desired data with the stuck-at resistance state, the number of readout errors in the system can be reduced. Previous works have used this property of stuck-at faults to mitigate their impact on the system [14].

According to the experiments using real fabricated ReRAM array in [33], stuck-at zeros (SA0) and stuck-at ones (SA1) can be clustered in the entire column/row or distributed randomly in a ReRAM array. This stochastic nature of SAF makes it hard to model using any specific distributions, thus, many previous studies [19], [23] have chosen the uniform distribution to model SAFs to reduce complexity during the fault estimation process. The same consideration is also applied in our paper. Furthermore, since the proposed method does not focus on any specific type of eNVMs or SAFs, using uniform distribution for evaluations of SAFs allows CRAFT to be generalized and applicable for any use case.

2) *Impact of Stuck-at Faults on Accuracy of DNNs*: A small number of stuck-at faults, especially in the critical bits, can cause a catastrophic change in DNN models accuracy [18], [23]. Fig. 3 shows the impact of stuck-at faults on different DNN models' accuracy evaluated on CIFAR-10 (a popular image dataset for computer vision). As illustrated in the figure, when the Bit Error Rate (BER) for stuck-at faults increases to a certain point, the classification error of the DNN model increases exponentially. For example, for ResNet-18 (a state-of-the-art Convolutional Neural Network) [3], when the BER increases beyond 2×10^{-6} , the classification error stays at around 90% level, which is the same as if the network is randomly guessing the results regardless of the input and trained parameters. Similar results are observed when considering different DNN models or datasets. Experimental details with additional results are discussed in Section V.

III. RELATED WORKS

Several memory-centric works have been proposed to address the problem of stuck-at faults in emerging NVMs. Error Correcting Pointers (ECP) [12] detect and correct stuck-at faults by keeping the stuck-at cell address and data in additional storage. [14] presents a method to enhance the correction capability of an ECC by a simple inversion operation. SAFER [13] dynamically partitions the data such that only single error is presented in each partition and then uses a single-bit error correction code for recovery. [34] proposes a method to reduce the storage overhead of ECP by allocating different number of error correction entries to different lines according to the number of hard errors presented

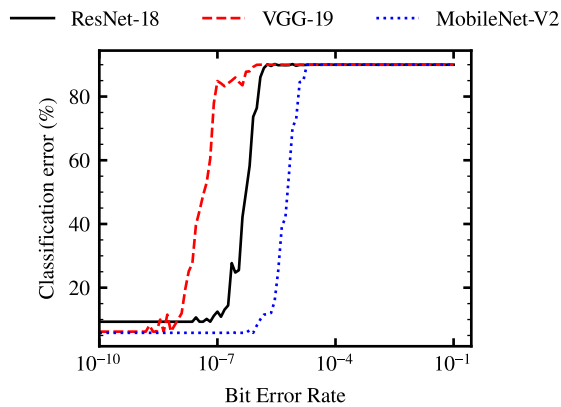


Fig. 3. Impact of Stuck-at faults on accuracy of different DNN models (ResNet-18, VGG-19 and MobileNet-V2 on CIFAR-10)

in the line. Although being effective in tolerating stuck-at faults of eNVMs-based system, these works often incur a large storage/hardware overhead, which is infeasible in the case of most resource-constrained edge devices. The proposed encoding techniques adds a minimum hardware overhead (1.17% in terms of storage) to the system yet still efficiently enhances the robustness of DNN against SAFs.

Apart from memory-centric approaches, several works have been proposed to enhance the stuck-at faults tolerance capability of DNN based systems. The work in [19] exploits the self-healing capability of DNNs and proposes a retraining method to reduce the impact of stuck-at cells in DNN accelerators. Such scheme can be effective in recovering the accuracy degradation from stuck-at errors; however, re-training is required when implementing these techniques, which is difficult to do when the DNN has been deployed to the edge devices. [35] redesigns the traditional error correction output code of DNNs using a collaborative logistic classifier, thus enhancing the DNN robustness against stuck-at faults. Despite being effective, this work also requires re-training (i.e., fine-tuning) to recover the accuracy impacted by SAFs. The need for re-training does not apply to the proposed fault-tolerant technique in this paper, since it is designed specifically for DNN inference on resource-constrained edge devices.

More relevant to our proposed techniques, many data remapping and redundancy based techniques have been proposed [20]–[22]. Specifically, [20] introduces mapping technique and redundant crossbar arrays to compensate the accuracy loss of DNN model caused by the stuck-at faults in ReRAM crossbar array. Since this technique utilizes redundant crossbar array that has the same size as the original array, the storage overhead and energy consumption of such technique is considerably large. [21] classifies weights according to their criticality to the model’s accuracy, remaps the significant weights to the fault-free memory cells and fine-tunes the DNN model to enhance the accuracy. By relying on the criticality of DNNs to address SAFs, such work is able to ease the re-training process and reduce storage overhead. Nonetheless, the storage overhead caused by such technique can still be as large as 5%, which is much higher compared to our proposed technique. The method in [23] uses matrix transformations to make the

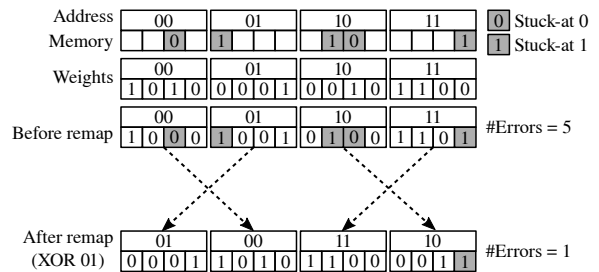


Fig. 4. Example of Intra-Block Address Remapping technique

weights in the ReRAM crossbar array more robust to stuck-at faults. Similar to other redundancy-based methods, [23] also adds a significant amount of hardware overhead compared to the proposed technique. Specifically, such a scheme can come at the expense of $8.19 \times$ power consumption and $9.23 \times$ area overhead. The proposed technique in this paper only adds six additional bits for encoding/decoding, making it highly efficient in term of storage/energy overhead

The existing memory-centric methods as well as DNN focused techniques often either require a large amount of additional hardware overhead or costly retraining process. In this paper, we propose a set of techniques that incur only minimal hardware overhead yet effectively enhance the fault-tolerance capability of DNNs in the presences of stuck-at faults. Moreover, the proposed techniques are orthogonal to the existing methods and can be implemented together to further enhance the robustness of DNNs.

IV. CRAFT: CRITICALITY-AWARE FAULT-TOLERANCE ENHANCEMENT TECHNIQUES

The state of a stuck-at cell can be detected (by a read operation) but cannot be re-programmed to a different state. Leveraging this fact, we present multiple remapping and encoding techniques, collectively named CRAFT, to reduce the number of stuck-at errors in the DNN parameters (weights and biases). The proposed fault-tolerance techniques include *Intra-Block Address Remapping*, *Weight Inversion* and *Criticality-Aware Bits Switching*.

A. Intra-Block Address Remapping

The parameters (weights and biases) of DNN are often stored as a group in a data block. For instance, a typical 64B cache line sized data block can store up to sixteen 32-bit floating point weights/biases. The proposed remapping method operates within the typical data block to preserve memory access locality. Fig. 4 shows a simple example of 2-bit address remapping using the proposed *Intra-Block Address Remapping* technique. An XOR operation of the address of each weight within a data block remaps the weight to a different location within the same block. This remapping allows to reduce the stuck-at faults by increasing the number of stuck-at cell states aligned with the desired bits. As shown in the figure, without remapping, the number of error due to stuck-at cells is five. After the proposed remapping technique (XOR the address with 01) is applied, the resulting errors due to stuck-at faults is reduced to one only. By considering multiple mappings

	W_0	W_1	W_2	W_3	W_4	W_5	W_6	W_7	W_8	W_9	W_{10}	W_{11}	W_{12}	W_{13}	W_{14}	W_{15}
Initial address	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111
XOR 0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000	0000
Remapped address	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111
	W_0	W_1	W_2	W_3	W_4	W_5	W_6	W_7	W_8	W_9	W_{10}	W_{11}	W_{12}	W_{13}	W_{14}	W_{15}
Initial address	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111
XOR 0001	0001	0001	0001	0001	0001	0001	0001	0001	0001	0001	0001	0001	0001	0001	0001	0001
Remapped address	0001	0000	0011	0010	0101	0100	0111	0110	1001	1000	1011	1010	1101	1100	1111	1110
	W_1	W_0	W_3	W_2	W_5	W_4	W_7	W_6	W_9	W_8	W_{11}	W_{10}	W_{13}	W_{12}	W_{15}	W_{14}
	⋮															
	W_0	W_1	W_2	W_3	W_4	W_5	W_6	W_7	W_8	W_9	W_{10}	W_{11}	W_{12}	W_{13}	W_{14}	W_{15}
Initial address	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111
XOR 1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111	1111
Remapped address	1111	1110	1101	1100	1011	1010	1001	1000	0111	0110	0101	0100	0011	0010	0001	0000
	W_{15}	W_{14}	W_{13}	W_{12}	W_{11}	W_{10}	W_9	W_8	W_7	W_6	W_5	W_4	W_3	W_2	W_1	W_0

Fig. 5. Intra-Block Address Remapping for 16 weights using 4-bit XOR operation

Weight	$\boxed{1}\boxed{0}\boxed{1}\boxed{1}\boxed{0}\boxed{0}\boxed{1}\boxed{1}\boxed{1}$	Value = 179
Three stuck-at faults	$\boxed{1}\boxed{0}\boxed{1}\boxed{1}\boxed{0}\boxed{1}\boxed{0}\boxed{0}$	Value = 180, Deviation = 1
One stuck-at fault	$\boxed{1}\boxed{0}\boxed{0}\boxed{1}\boxed{0}\boxed{0}\boxed{1}\boxed{1}$	Value = 147, Deviation = 32
	$\boxed{0}$ Stuck-at 0 $\boxed{1}$ Stuck-at 1	

Fig. 6. Example of fault significance in a data block. Three errors in insignificant bit positions result in insignificant change in value while one error in significant bit position causes large deviation from the actual value.

through different XOR operations, the number of errors can be further reduced. The proposed method finally chooses the mapping which causes minimal impact (instead of the fewest errors, as explained in the next section) on DNNs accuracy. Fig. 5 illustrates the proposed remapping technique using 4-bit XOR operation. As shown in the figure, sixteen different mappings can be obtained using 4-bit XOR operation.

1) Minimizing the Impact of Faults on DNN's Accuracy:

Errors that occur in significant bits of DNN parameters are more harmful to the network accuracy than the errors in insignificant bits. This can be easily understood using a simple example shown in Fig. 6. As illustrated in the figure, frequent stuck-at faults in the insignificant bit positions result in only a small change from the actual weight value while fewer faults in the significant bit positions cause a greater change to the weight value. This is indicated by measuring the deviation in decimal value of the erroneous weight from that of the actual weight. As shown in Fig. 6, three stuck-at faults in the insignificant bit position result in deviation of 1 (in decimal), while a single stuck-at fault in the significant bit position causes a deviation of 32 (in decimal). In floating-point representation, this criticality difference is even more evident due to the greater difference in the significance of exponent bits as compared to the other bit positions. Furthermore, in case of DNN, certain weights are more critical than the other weights, thus, merely minimizing the number of faults in the memory would not be always helpful. Therefore, the proposed remapping method chooses to minimize the deviation (from the original value) of weights instead of simply minimizing the number of stuck-at faults presented in memory. The proposed deviation minimization technique can be formulated as:

Address	00	01	10	11	$\boxed{0}$ Stuck-at 0
Memory	$\boxed{0}$	$\boxed{1}$	$\boxed{1}\boxed{0}$	$\boxed{1}$	$\boxed{1}$ Stuck-at 1
Weights	W_1	W_2	W_3	W_4	
	$0\ 1\ 1\ 1$	$0\ 1\ 0\ 1$	$1\ 0\ 1\ 1$	$0\ 1\ 1\ 0$	
Initial	W_1	W_2	W_3	W_4	Errors = 5
	$0\ 1\ 0\ 1$	$1\ 1\ 0\ 1$	$1\ 1\ 0\ 1$	$0\ 1\ 1\ 1$	Deviation = 2+8+4+2+1=17
Address	W_2	W_1	W_4	W_3	Errors = 2
XOR 01	$0\ 1\ 0\ 1$	$1\ 1\ 1\ 1$	$0\ 1\ 0\ 0$	$1\ 0\ 0\ 1$	Deviation = 8+2 = 10
Address	W_3	W_4	W_1	W_2	Errors = 3
XOR 10	$1\ 0\ 0\ 1$	$1\ 1\ 1\ 0$	$0\ 1\ 0\ 1$	$0\ 1\ 0\ 1$	Deviation = 2+8+2 = 12
Address	W_4	W_3	W_2	W_1	Errors = 1
XOR 11	$0\ 1\ 0\ 0$	$1\ 0\ 1\ 1$	$0\ 1\ 0\ 1$	$0\ 1\ 1\ 1$	Deviation = 2
	↓				Minimal mapping

Fig. 7. An example of Intra-Block Address Remapping technique with deviation minimization. Net deviation for each mapping is calculated by summing up all weight value deviations in the data block. Mapping with minimum net deviation is chosen as minimal mapping for inference.

$$w' \leftarrow w_{r_i} \text{ s.t } \delta = \min_{\delta \in \Delta} \sum_{i=0}^N |w_{r_i} - w_{o_i}| \quad (1)$$

where w' is the final weight that is used for inference, w_{r_i} and w_{o_i} refer to the new weight after remapping (obtained while considering stuck-at faults) and the original weight, respectively. N is the number of weights within a selected data block. δ is the minimum net deviation and Δ is the set of all possible net deviations for different remappings of the weights.

The proposed Intra-Block Address Remapping technique with a deviation minimization algorithm is illustrated in Fig. 7. For illustration simplicity, four 4-bit weights with 2-bit XOR operation for address remapping are depicted in the example. Five stuck-at faults (three SA1 and two SA0) are randomly distributed across the memory block, as shown in the figure. Initially, without any remapping, the readout data results in five errors which leads to the net deviation of 13 (in decimal). Using 2-bit XOR operation, four possible mappings can be obtained. As illustrated in the figure, the mapping which uses XOR 11 operation leads to the minimum net deviation (= 2 in decimal) from the actual weight. Therefore, the proposed

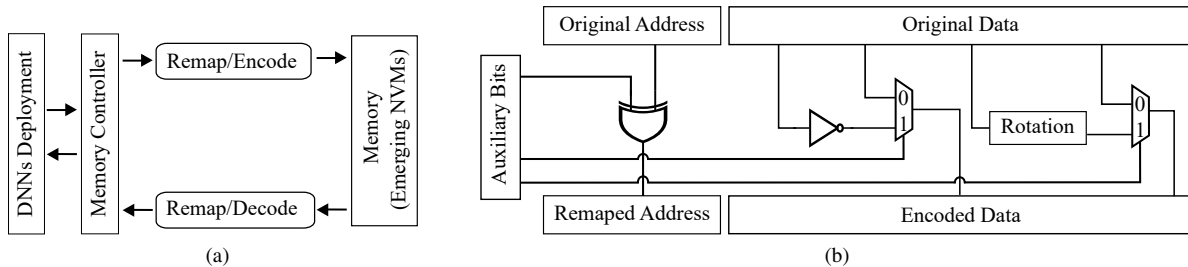


Fig. 11. Implementation of CRAFT (a) Overall Architecture (b) Remapping and Encoding Logic

TABLE I
SPECIFICATIONS OF EVALUATED DNNs AND DATASETS

Dataset	Model	Input Dimension	Parameters size (MB)	Classification Error (%)	Precision
CIFAR-10	ResNet-18	3x32x32	43	9.31	32-bit floating-point
	VGG-19	3x32x32	533	6.19	32-bit floating-point
	MobileNet-V2	3x32x32	9	5.83	32-bit floating-point
Imagenet	ResNet-50	3x224x224	25.6	24.02	8-bit quantized
	VGG-19	3x224x224	143.72	27.64	8-bit quantized
	Inception-V4	3x224x224	43.03	20.06	8-bit quantized

A. Experimental Setup

The simulations for stuck-at faults in DNNs are performed by using the Pytorch framework [37]. All simulations are performed on an Intel® Xeon® CPU E5-1650 v4 with two Nvidia Titan XP GPUs having 24Gb of RAM. Table I lists the specifications of the evaluated DNN models and datasets. The proposed fault-tolerance enhancement techniques of the CRAFT architecture are evaluated for various state-of-the-art DNNs using popular datasets. The DNN models used in the experiments include *ResNet*, *VGG*, *MobileNet* and *Inception*. The two datasets considered are *CIFAR-10* and *ImageNet*. The *CIFAR-10* dataset contains 60,000 RGB images, of which, 50,000 images are for training and 10,000 images are for testing. Training set images are preprocessed by padding 4 pixels along the height and width and randomly cropping to a 32×32 patch. Furthermore, random horizontal flip operation is also performed on the training set images. DNN models evaluated on *CIFAR-10* dataset (i.e., ResNet-18, VGG-19 and MobileNet-V2) are trained using stochastic gradient descent with 0.9 momentum. The cross entropy loss function is used as the objective function to classify ten classes of the input images [3]. The learning rate is kept constant at 0.1 during the training process.

Beside the *CIFAR-10* dataset, the proposed method is also evaluated on the *ImageNet* dataset, which consists of 1.3M images in the training set and 50,000 images in the test set. Training images are first randomly cropped to a 224×224 patch and then followed by a random horizontal flip operation. Both training set and test set images are normalized using channel wise normalization for zero mean and unit standard deviation. During the training process, the network is learned by using the cross entropy loss as the objective function and minimized by the stochastic gradient descent algorithms with a momentum of 0.9. The learning rate is also kept at 0.1 during training on the *ImageNet* dataset. After being trained, the network parameters are quantized to a lower precision (8-

bit quantization) from 32-bit floating point representation. The baseline classification error of the evaluated models are shown in Table. I.

Quantization is a widespread technique, considered in many practical implementations, to reduce the computational overhead in DNNs. The network parameters are quantized using fewer bits resulting in a reduced precision. To evaluate the proposed method on quantized networks, we implement an 8-bit quantization scheme on DNNs that are evaluated on *ImageNet* dataset. For quantized DNN models, a simple yet effective quantization scheme introduced in [26] is considered. Using this scheme, we implement 8-bit, weight only post-training quantization with asymmetric and per-layer mode.

B. Stuck-at Faults Simulation Framework

Stuck-at faults in emerging memory-based DNNs are simulated using a customized fault-injection framework built upon the validated Ares framework introduced in [38]. The experiments mainly consist of two stages. The first stage constructs various DNNs for experimentation using different weight precision and obtains possible parameter remappings/encodings based on the proposed methods presented in Sec. IV.

In the second stage, the Ares framework is applied during DNN inference with the consideration of stuck-at faults in emerging NVMs. The stuck-at faults are assumed to be uniformly distributed in the network parameters, as mentioned in Sec. II-C1. The proposed techniques are applied by considering the remapping/encoding configuration that leads to minimum net deviation in stage 1. For each bit error rate, the network classification error is averaged over 100 fault-injection simulations.

C. Comparison with the Error Correcting Techniques

The effectiveness of the proposed techniques is compared with the relevant existing Error Correcting Coding (ECC)

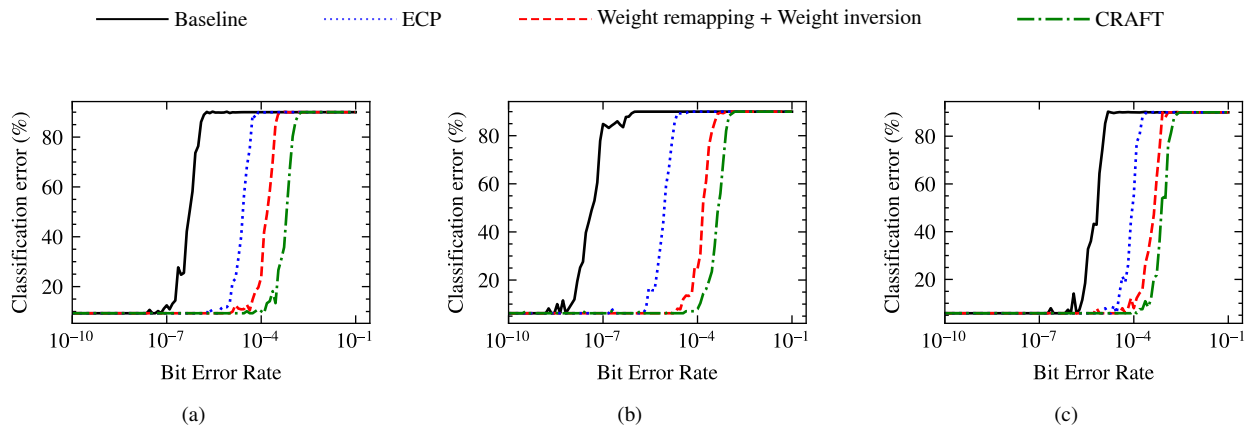


Fig. 12. Comparison of different fault-tolerance techniques for various DNNs using CIFAR-10 dataset and 32-bit floating-point parameters (a) ResNet-18 (b) VGG-19 (c) MobileNet-V2

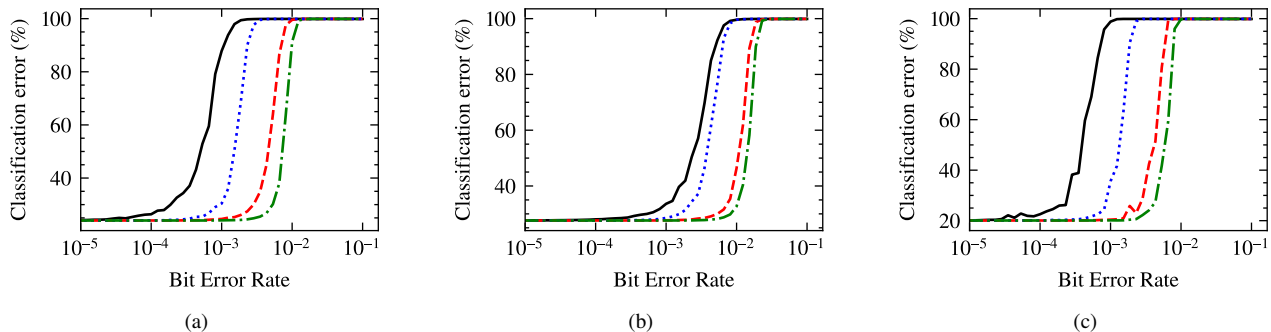


Fig. 13. Comparison of different fault-tolerance techniques for various DNNs using Imagenet dataset and 8-bit quantized parameters (a) ResNet-50 (b) VGG-19 (c) Inception-V4

solutions for the same purpose. As discussed earlier, the proposed techniques are flexible in terms of trade-off between the storage overhead and degree of robustness against stuck-at faults (more auxiliary bit can be used to enhance the fault-tolerance capability of the system at the cost of increased storage overhead). Therefore, we compare the proposed method with existing ECC techniques with similar storage overhead.

In conventional memory system, ECC code such as (72, 64) Hamming code is often used for addressing soft errors. A typical (72, 64) Hamming code normally incurs more than 10% of storage overhead, which is $10\times$ larger than the combined overhead of all of the proposed techniques of CRAFT. Moreover, ECC schemes for addressing soft errors in conventional memory system are not as effective in emerging NVMs, as shown and discussed in [12]. Instead, Error Correcting Pointers (ECP) are often preferred for tolerating hard errors in emerging NVMs. For a fair comparison, we benchmark the proposed methods with the ECP variant that incurs a comparable storage overhead. For d -bit data, ECP_n is able to correct up to n bits with the storage overhead of $\frac{1+n+n \cdot \lceil \log_2 d \rceil}{d}$. For our experiments, we consider ECP_1 which can correct a single bit error in 512-bit data block at the cost of 2.15% storage overhead, which is still higher than that of CRAFT.

As discussed in Sec. IV, for 32-bit data (one full-precision weight or four quantized weights), the proposed techniques together require approximately 1.17% storage overhead. This

amount is $2\times$ less than the evaluated ECP. Despite having a minimal storage overhead, CRAFT is found to be more effective as compared to ECP, regardless of the DNN model type or dataset size. The detailed results are presented in the next section.

D. Results and Discussion

The evaluation results of the proposed methods for CIFAR-10 and ImageNet are shown in Fig. 12 and Fig. 13, respectively. The solid black line shows the baseline models (considering stuck-at faults) without any error-correction scheme. The dotted blue line indicates the models which incorporate ECP to mitigate stuck-at faults. The red dashed line illustrates models that use *Intra-Block Addressing Remapping* and *Weight Inversion* technique for stuck-at fault-tolerance. The green dash-dotted line shows results for CRAFT architecture which includes *Weight Remapping + Weight Inversion + Criticality-Aware Bits Switching*. As seen from the results, CRAFT improves the robustness of the baseline models significantly and outperforms the ECP method by a significant margin.

Specifically, Fig. 12(a) shows the evaluation results of different fault-tolerance methods for ResNet-18 using CIFAR-10 dataset. The baseline and ECP classification error starts to increase exponentially at the bit error rate (BER) of 1.5×10^{-7} and 5×10^{-6} , respectively. On the other hand, the classification error can be maintained below 10% at around 5×10^{-5} BER when Weight Remapping and Weight Inversion are applied and

TABLE II
SUMMARY OF ROBUSTNESS ENHANCEMENT BY DIFFERENT FAULT-TOLERANCE TECHNIQUES FOR VARIOUS DNNs AND DATASETS

Dataset	Model	Robustness Improvement			
		Baseline	ECP	Weight Remapping + Weight Inversion [24]	CRAFT
CIFAR-10	ResNet-18	1x	81x	351x	1233x
	VGG-19	1x	284x	3511x	12320x
	MobileNet-V2	1x	23x	53x	231x
ImageNet	ResNet-50	1x	4x	15x	29x
	VGG-19	1x	2x	8x	12x
	Inception-V4	1x	5x	10x	23x

at 2×10^{-4} when used together with Criticality-Aware Bits Switching in CRAFT. In other words, CRAFT can increase the fault-tolerance by up to more than 1200 \times compared to the baseline and 3 \times compared to only using Weight Remapping + Inversion methods. A similar trend can also be observed in other networks using the CIFAR-10 dataset such as VGG-19 and MobileNet-V2. For example, in VGG-19, while ECP can only improve the robustness of the model up to 284 \times , CRAFT can increase the robustness to 12,320 \times compared with the baseline model, which is orders of magnitude higher than ECP. The efficiency of CRAFT compared to the baseline model when using MobileNet-V2 is found to be 231 \times , which is 10 \times higher than ECP and 4 \times better than previously proposed methods [24].

To confirm the general applicability of the proposed techniques, we also perform different experiments on larger DNNs (ResNet-50, Inception-V4, etc...) with larger dataset (i.e., ImageNet) (Fig. 13). As discussed in Sec. IV-C1, the 8-bit quantized networks show a significant improvement in robustness compared to full-precision networks in Fig. 12 due to the smaller dynamic range representation. Regardless of such property, the proposed techniques still ensure a significant increase in robustness against stuck-at faults, as shown in Fig. 13. For example, for ResNet-50 using ImageNet (Fig. 13(a)), the network accuracy starts to drop at around 2×10^{-4} BER while CRAFT can maintain the accuracy up to 5×10^{-3} BER, indicating 29 \times more robustness than the baseline model. This trend is found to be consistent in other DNNs evaluated on ImageNet dataset. Specifically, CRAFT can enhance the robustness of VGG-19 and Inception-V4 up to 12 \times and 23 \times , respectively. A summary of the robustness improvement over the baseline configurations using different fault-tolerance techniques is given in Table. II. The improvement in robustness is obtained by measuring the BER at which the technique causes the classification error to increase by more than 5%. As shown in the table, CRAFT outperforms the existing techniques by orders of magnitude (up to 10^4 times over the baseline). While achieving a significant improvement in robustness against SAFs, as mentioned in Sec. IV-D, CRAFT incurs a minimum amount of storage overhead ($\approx 1.17\%$), which makes CRAFT highly practical when implementing in resource-constrained edge devices.

Stuck-at faults in emerging NVMs are expected to become more frequent as technology scaling to smaller nodes. The two distinguishing features of CRAFT are consideration for the criticality of errors and flexibility to address more errors

by employing a fine-grained remapping using smaller block size. This approach makes CRAFT a robust and scalable method which can tackle future trends of hard errors in NVMs. Finally, the proposed techniques of CRAFT are orthogonal to the existing error correcting techniques, therefore, further enhancement in the fault-tolerance can be achieved by implementing CRAFT along with the conventional techniques.

VI. CONCLUSION

Hard errors such as stuck-at faults can severely impact the accuracy of Deep Neural Networks based systems which use emerging Non-Volatile Memories. This paper introduces a set of robust techniques, collectively named CRAFT, to enhance the error-tolerability of DNNs against stuck-at faults. The proposed techniques are simple and light-weight yet effective to tackle the problem of stuck-at faults in DNNs-based system. Working in a hierarchical manner, the proposed CRAFT architecture remaps the weights, encodes them using a simple inversion method and switches the bits within the weights based on their criticality, thus minimizing the impact of stuck-at faults on neural network's accuracy. The evaluation results show that CRAFT is able to enhance the robustness of the system by orders of magnitude. Specifically, for DNNs evaluated on CIFAR-10, CRAFT can increase the robustness by up to 10^4 times compared to the baseline model. For DNNs using ImageNet, CRAFT enhances the robustness of the model by up to 29 times. Being orthogonal, the proposed techniques of CRAFT can be easily incorporated with other existing methods to further increase the fault-tolerance of DNNs.

REFERENCES

- [1] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," 2015.
- [2] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, pp. 525–542.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [4] S. Jeong, S. Kang, and J.-S. Yang, "Pair: Pin-aligned in-dram ecc architecture using expandability of reed-solomon code," in *2020 57th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2020, pp. 1–6.
- [5] S. Yu and P. Chen, "Emerging memory technologies: Recent trends and prospects," *IEEE Solid-State Circuits Magazine*, vol. 8, no. 2, pp. 43–56, 2016.

- [6] B. Rajendran and F. Alibart, "Neuromorphic computing based on emerging memory technologies," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 6, no. 2, pp. 198–211, 2016.
- [7] A. Shafiee, A. Nag, N. Muralimanohar, R. Balasubramonian, J. P. Strachan, M. Hu, R. S. Williams, and V. Srikumar, "Isaac: A convolutional neural network accelerator with in-situ analog arithmetic in crossbars," *ACM SIGARCH Computer Architecture News*, vol. 44, no. 3, pp. 14–26, 2016.
- [8] V. Joshi, M. Le Gallo, S. Haefeli, I. Boybat, S. R. Nandakumar, C. Piveteau, M. Dazzi, B. Rajendran, A. Sebastian, and E. Eleftheriou, "Accurate deep neural network inference using computational phase-change memory," *Nature communications*, vol. 11, no. 1, pp. 1–13, 2020.
- [9] R. F. Freitas and W. W. Wilcke, "Storage-class memory: The next storage system technology," *IBM Journal of Research and Development*, vol. 52, no. 4.5, pp. 439–447, 2008.
- [10] D. Ielmini and H.-S. P. Wong, "In-memory computing with resistive switching devices," *Nature Electronics*, vol. 1, no. 6, pp. 333–343, 2018.
- [11] T. Kwon, M. Imran, and J.-S. Yang, "Reliability enhanced heterogeneous phase change memory architecture for performance and energy efficiency," *IEEE Transactions on Computers*, vol. 70, no. 9, pp. 1388–1400, sep 2021. [Online]. Available: <https://doi.org/10.1109/2Ftc.2020.3009498>
- [12] S. Schechter, G. H. Loh, K. Strauss, and D. Burger, "Use ecp, not ecc, for hard failures in resistive memories," in *Proceedings of the 37th Annual International Symposium on Computer Architecture*, ser. ISCA' 10. New York, NY, USA: Association for Computing Machinery, 2010, pp. 141–152.
- [13] N. H. Seong, D. H. Woo, V. Srinivasan, J. A. Rivers, and H. S. Lee, "Safer: Stuck-at-fault error recovery for memories," in *2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture*, 2010, pp. 115–124.
- [14] R. Maddah, S. Cho, and R. Melhem, "Power of one bit: Increasing error correction capability with data inversion," in *2013 IEEE 19th Pacific Rim International Symposium on Dependable Computing*, 2013, pp. 216–225.
- [15] R. Maddah, R. Melhem, and S. Cho, "Rdis: Tolerating many stuck-at faults in resistive memory," *IEEE Transactions on Computers*, vol. 64, no. 3, pp. 847–861, 2015.
- [16] M. Asadinia, M. Jalili, and H. Sarbazi-Azad, "Bless: A simple and efficient scheme for prolonging pcm lifetime," in *Proceedings of the 53rd Annual Design Automation Conference*, ser. DAC '16. New York, NY, USA: Association for Computing Machinery, 2016. [Online]. Available: <https://doi.org/10.1145/2897937.2897993>
- [17] S. Mittal, "A survey on modeling and improving reliability of dnn algorithms and accelerators," *Journal of Systems Architecture*, vol. 104, p. 101689, 2020.
- [18] G. Li, S. K. S. Hari, M. Sullivan, T. Tsai, K. Pattabiraman, J. Emer, and S. W. Keckler, "Understanding error propagation in deep learning neural network (dnn) accelerators and applications," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 2017, pp. 1–12.
- [19] L. Chen, J. Li, Y. Chen, Q. Deng, J. Shen, X. Liang, and L. Jiang, "Accelerator-friendly neural-network training: Learning variations and defects in rram crossbar," in *Design, Automation Test in Europe Conference Exhibition (DATE), 2017*, 2017, pp. 19–24.
- [20] W. Huangfu, L. Xia, M. Cheng, X. Yin, T. Tang, B. Li, K. Chakrabarty, Y. Xie, Y. Wang, and H. Yang, "Computation-oriented fault-tolerance schemes for rram computing systems," in *2017 22nd Asia and South Pacific Design Automation Conference (ASP-DAC)*, 2017, pp. 794–799.
- [21] C. Liu, M. Hu, J. P. Strachan, and H. Li, "Rescuing memristor-based neuromorphic design with high defects," in *2017 54th ACM/EDAC/IEEE Design Automation Conference (DAC)*, 2017, pp. 1–6.
- [22] L. Xia, W. Huangfu, T. Tang, X. Yin, K. Chakrabarty, Y. Xie, Y. Wang, and H. Yang, "Stuck-at fault tolerance in rram computing systems," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 8, no. 1, pp. 102–115, 2018.
- [23] B. Zhang, N. Uysal, D. Fan, and R. Ewetz, "Handling stuck-at-faults in memristor crossbar arrays using matrix transformations," in *Proceedings of the 24th Asia and South Pacific Design Automation Conference*, ser. ASPDAC '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 438–443. [Online]. Available: <https://doi.org/10.1145/3287624.3287707>
- [24] T.-H. Nguyen, M. Imran, J. Choi, and J.-S. Yang, "Low-cost and effective fault-tolerance enhancement techniques for emerging memories-based deep neural networks," in *2021 58th ACM/IEEE Design Automation Conference (DAC)*. IEEE, dec 2021. [Online]. Available: <https://doi.org/10.1109/dac18074.2021.9586112>
- [25] D. Han, J. Kim, and J. Kim, "Deep pyramidal residual networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5927–5935.
- [26] R. Krishnamoorthi, "Quantizing deep convolutional networks for efficient inference: A whitepaper," *arXiv preprint arXiv:1806.08342*, 2018.
- [27] F. Pan, S. Gao, C. Chen, C. Song, and F. Zeng, "Recent progress in resistive random access memories: Materials, switching mechanisms, and performance," *Materials Science and Engineering: R: Reports*, vol. 83, pp. 1–59, 2014. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0927796X14000692>
- [28] M. Imran, T. Kwon, J. M. You, and J. S. Yang, "Flipcy: Efficient pattern redistribution for enhancing mlc pcm reliability and storage density," in *2019 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, 2019, pp. 1–7.
- [29] T. Kwon, M. Imran, and J.-S. Yang, "Cost-effective reliable mlc pcm architecture using virtual data based error correction," *IEEE Access*, vol. 8, pp. 44 006–44 018, 2020.
- [30] T.-H. Nguyen, M. Imran, and J.-S. Yang, "DynaPAT," in *Proceedings of the 41st IEEE/ACM International Conference on Computer-Aided Design*. ACM, oct 2022. [Online]. Available: <https://doi.org/10.1145/2F3508352.3549400>
- [31] T. Kwon, M. Imran, and J.-S. Yang, "Pattern-aware encoding for mlc pcm storage density, energy efficiency, and performance enhancement," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 39, no. 9, pp. 1855–1865, 2020.
- [32] M. Imran, T. Kwon, N. A. Toubia, and J.-S. Yang, "Cent: An efficient architecture to eliminate intra-array write disturbance in pcm," *IEEE Transactions on Computers*, pp. 1–1, 2021.
- [33] C.-Y. Chen, H.-C. Shih, C.-W. Wu, C.-H. Lin, P.-F. Chiu, S.-S. Sheu, and F. T. Chen, "Rram defect modeling and failure analysis based on march test and a novel squeeze-search scheme," *IEEE Transactions on Computers*, vol. 64, no. 1, pp. 180–190, 2014.
- [34] M. K. Qureshi, "Pay-as-you-go: Low-overhead hard-error correction for phase change memories," in *Proceedings of the 44th Annual IEEE/ACM International Symposium on Microarchitecture*, 2011, pp. 318–328.
- [35] T. Liu, W. Wen, L. Jiang, Y. Wang, C. Yang, and G. Quan, "A fault-tolerant neural network architecture," in *2019 56th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2019, pp. 1–6.
- [36] J.-S. Kim and J.-S. Yang, "Dris-3: Deep neural network reliability improvement scheme in 3d die-stacked memory based on fault analysis," in *2019 56th ACM/IEEE Design Automation Conference (DAC)*, 2019, pp. 1–6.
- [37] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.
- [38] B. Reagen, U. Gupta, L. Pentecost, P. Whatmough, S. K. Lee, N. Mulholland, D. Brooks, and G.-Y. Wei, "Ares: A framework for quantifying the resilience of deep neural networks," in *2018 55th ACM/ESDA/IEEE Design Automation Conference (DAC)*. IEEE, 2018, pp. 1–6.



processing-in-memory and computationally efficient implementations of deep learning algorithms.

Thai-Hoang Nguyen (Student Member, IEEE) received the BS degree in electrical engineering from Danang University of Science and Technology, the University of Danang, Vietnam in 2019. He joined Sungkyunkwan University, Suwon, Korea, in 2019 where he is currently working toward his PhD degree in electrical and computer engineering. He is the recipient of STEM scholarship for international graduate student by Sungkyunkwan University, from 2019. His current research interests include designing robust architecture for emerging memory technologies,



Muhammad Imran received the BS degree in electrical engineering from the University of Engineering and Technology, Lahore, Pakistan, in 2012 and the PhD degree in electronic and electrical engineering from Sungkyunkwan University, Suwon, South Korea, in 2020. He is currently an assistant professor at the National University of Sciences and Technology, Islamabad, Pakistan. He is recipient of the Superior Research Award at Sungkyunkwan University and Scholarship by the Higher Education Commission of Pakistan for his MS and PhD studies.

His current research interests include vector processor design and reliable computing architectures for deep learning.



Jaehyuk Choi Jaehyuk Choi (Associate Member, IEEE) received the B.S. degree in electrical engineering from Yonsei University, Seoul, South Korea, in 2004, the M.S. degree in electrical engineering and computer science from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 2006, the M.S. degree in electrical and computer engineering from the University of Minnesota, Minneapolis, MN, USA, in 2008, and the Ph.D. degree from the University of Michigan, Ann Arbor, MI, USA, in 2013. From 2013

to 2015, he was a Research Staff Member with the Samsung Advanced Institute of Technology (SAIT), Samsung Electronics, Suwon, South Korea, where he engaged in the development of depth sensors and low-power image sensors. In 2015, he joined the Department of Semiconductor Systems Engineering, Sungkyunkwan University, Suwon, as an Assistant Professor. His research interests include low-power circuits, CMOS sensors, and mixed-signal integrated circuits.



Joon-Sung Yang Joon-Sung Yang (Senior Member, IEEE) received the B.S. degree from Yonsei University, Seoul, South Korea, in 2003, and the M.S. and Ph.D. degrees from The University of Texas at Austin, Austin, TX, USA, in 2007 and 2009, respectively, all in electrical and computer engineering. After graduation, he worked at Intel Corporation, Austin, for four years. He was with Sungkyunwan University. He is currently an Associate Professor with Yonsei University. His research interests include memory architectures and efficient

deep learning architecture development. He was a recipient of the Korea Science and Engineering Foundation (KOSEF) Scholarship, in 2005. He received the Best Paper Award from the 2008 IEEE International Symposium on Defect and Fault Tolerance in VLSI systems and the 2016 IEEE International SoC Design Conference. He was nominated for the Best Paper Award from the 2013 IEEE VLSI Test Symposium.