# Exploration of Activation Fault Reliability in Quantized Systolic Array-Based DNN Accelerators

Mahdi Taheri[1], Natalia Cherezova[1], Mohammad Saeed Ansari[2], Maksim Jenihhin[1],
Ali Mahani[3,4], Masoud Daneshtalab[1,5], Jaan Raik[1]

[1]Tallinn University of Technology, Tallinn, Estonia
[2]University of Alberta, Edmonton, Canada,
[3]Shahid Bahonar University of Kerman, Kerman, Iran
[4]York University, Toronto, Canada
[5]Mälardalen University, Västerås, Sweden
[1]mahdi.taheri@taltech.ee

*Abstract*—The stringent requirements for the Deep Neural Networks (DNNs) accelerator's reliability stand along with the need for reducing the computational burden on the hardware platforms, i.e. reducing the energy consumption and execution time as well as increasing the efficiency of DNN accelerators. Moreover, the growing demand for specialized DNN accelerators with tailored requirements, particularly for safety-critical applications, necessitates a comprehensive design space exploration to enable the development of efficient and robust accelerators that meet those requirements. Therefore, the trade-off between hardware performance, i.e. area and delay, and the reliability of the DNN accelerator implementation becomes critical and requires tools for analysis. This paper presents a comprehensive methodology for exploring and enabling a holistic assessment of the trilateral impact of quantization on model accuracy, activation fault reliability, and hardware efficiency. A fully automated framework is introduced that is capable of applying various quantization-aware techniques, fault injection, and hardware implementation, thus enabling the measurement of hardware parameters. Moreover, this paper proposes a novel lightweight protection technique integrated within the framework to ensure the dependable deployment of the final systolic-array-based FPGA implementation. The experiments on established benchmarks demonstrate the analysis flow and the profound implications of quantization on reliability, hardware performance, and network accuracy, particularly concerning the transient faults in the network's activations.

*Index Terms*—deep neural networks, design space exploration, quantization, fault simulation, reliability assessment

## I. INTRODUCTION

In the past decades, Deep Neural Networks (DNNs) demonstrated a significant improvement in accuracy by adopting intense parameterized models [1]. As a consequence, the size of these models has drastically increased, imposing challenges in deploying them on resource-constrained platforms [2]. FPGAs are a widely used solution for flexible and efficient DNN accelerator implementations and have shown superior hardware performance in terms of latency and power [3]. In practice, deployment of an FPGA-based DNN accelerator for the safety- and mission-critical applications (e.g., autonomous driving) requires addressing the trade-off between different
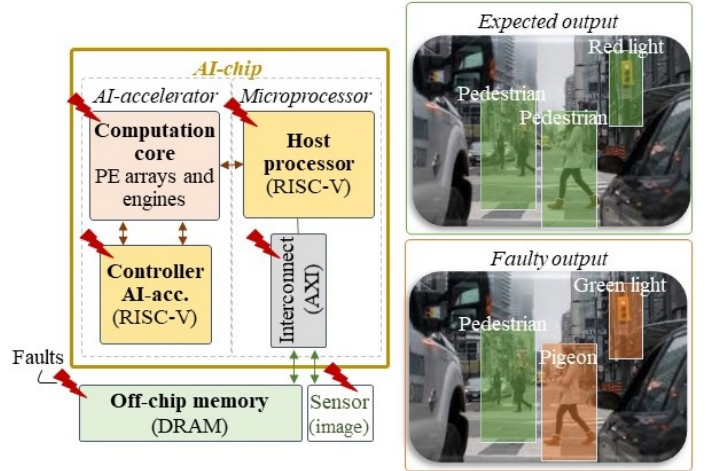


Fig. 1: Hardware-induced reliability threats in an example DNN accelerator and their possible impact on the output

design parameters of *hardware performance*, e.g., area, power, delay, and *reliability*. A compromise between conflicting requirements can be achieved by simplifying the implementation to sacrifice the precision of results but benefiting from lower resource utilization, energy consumption, and higher system efficiency. *Quantization* is one such concept that is being widely used in neural network deployments [4]. Quantization is used to compress the model for storage and computation reduction. However, recent research shows that faults in memory can cause a significant drop in DNN accuracy, which raises concern about the impact of quantization on the reliability of the network [5].

The reliability of DNN accelerators expresses their ability to produce correct outputs in the presence of hardware faults originating from various phenomena, e.g., radiation-induced soft errors in memory or logic [6]. DNNs are known to be inherently fault-resilient due to the high number of learning process iterations and several parallel neurons with multiple

computation units. Nevertheless, faults may impact the output accuracy of DNNs drastically [7], and in the case of resource-constrained critical applications, the reliability of DNNs is required to be evaluated and guaranteed [8]. The complexity of such evaluation motivates an *automated toolchain* with quantization and reliability analysis to support *Design Space Exploration (DSE)* for DNN accelerators already at the early design stage, i.e. starting from a high-level description, followed by providing an FPGA prototype for the selected design.

While the protection of weights stored in ROM can be ensured through error correction codes (ECC) or similar protection techniques, the dynamic nature of activations, which are stored for a short period of time in usually unprotected memories, poses a critical concern. Thus, it is crucial to thoroughly investigate the consequences of faults in the network's activations.

This paper presents a framework containing a fully automated toolchain to perform a study on the impact of quantization on network accuracy, hardware performance, and reliability drop in the presence of activation faults (Fig. 1) in systolic-array-based FPGA accelerators. To the best of our knowledge, this is the first framework that holistically considers those parameters. A novel lightweight mitigation technique is proposed and integrated into the framework to study potential trade-offs of compensating the reliability drops. The proposed methodology enables the analysis both at the level of the network model and at the level of individual layers of the network.

This framework is empowered by techniques for quantizing the networks and restricting the activation ranges to be limited to a certain level throughout the whole network execution by applying an extra scaling function in the network inference. This framework uses the high-level description of a DNN as an input and is capable of providing a transient-fault-resilient systolic-array-based FPGA implementation of the network utilizing the design parameters selected by the DSE. The main contributions in this work are as follows:

- A methodology for holistic exploration of quantization and reliability trade-offs in systolic-array implementation that enables assessing the trilateral impact of quantization on accuracy, activation fault reliability, and hardware performance.
- A fully-automated framework that is capable of applying quantization-aware training, post-training quantization, range-restriction, fault simulation, and implementing the whole methodology down to hardware implementation to measure actual hardware parameters like area, latency, etc.
- A lightweight and effective protection technique is developed and adopted in the framework toolchain to provide the final reliable systolic-array-based FPGA implementation of the network
- Demonstration and analysis of the results on the impact of quantization on reliability, hardware performance, and accuracy of the neural networks due to the transient faults in the activations for two well-known benchmarks.

The rest of the paper is organized as follows. Related works are discussed in Section II, the methodology and framework are presented in Section III, the experimental setup and results are provided in Section IV, and finally, the work is concluded in Section V.

## II. RELATED WORKS

### A. DNN reliability and quantization studies

Several works examine the impact of different fault models on the basis of a number of layers in DNNs and different data types [9]. Investigation into the effects of data precision is done in [10], where authors conducted a comparison of the resilience of FP16, FP32, and FP64 in the context of Matrix Multiplication. Their findings indicated that the reduction of precision not only enhances GPU performance and efficiency but also contributes to its overall resilience.

Another study [11] involved the deployment of MNIST CNN on FPGAs utilizing FP32, FP16. The results of the experiment demonstrated that decreasing the data precision in CNNs can lead to a substantial enhancement in overall resilience. This improvement was attributed to the reduced memory usage. Furthermore, [12] noted that the application of binary quantization to weights in convolutional layers results in decreased vulnerability factors, although it does increase the criticality of faults. [13] showed that the impact of faults is higher in most significant bits (MSBs) and with more aggressive compression the most significant bits are more probable to be exposed to faults. The aforementioned works show that quantization from higher data representations like FP32 down to INT16 has a positive impact on the performance and overall resilience, though *on the lower quantization ranks this matter should be studied and is not always impacting positively on the resilience*

In [14], it is shown that in some cases, the impact of the faults in the weight memories of a DNN can be negligible. Even though in the above-mentioned works, *impact of faults (soft errors modeled as bit flips) in the weights of a DNN during inference is examined*, to further enhance our comprehension of the impact of quantization on the reliability of DNNs in systolic-array-based DNN accelerators, this work is enriched with an FI engine capable of injecting faults into the activations of the DNNs in the systolic architecture.

### B. Fault mitigation techniques

The process of quantization and outlier regularization offers the potential to restrict the numerical range within a DNN, thereby eliminating the possibility of generating excessively large values due to faults [5].

Hoang et al. analyzed how various boundary values affect the network's accuracy. They have found that the best boundary values for each layer are not necessarily the maximum values of the layers' activations [15]. Hence, they propose an interval search algorithm to find appropriate boundary values for the ReLU activation function at each layer, named FT-ClipAct. The proposed clipped activation function maps their outputs to 0 if activations exceed the boundaries. Although
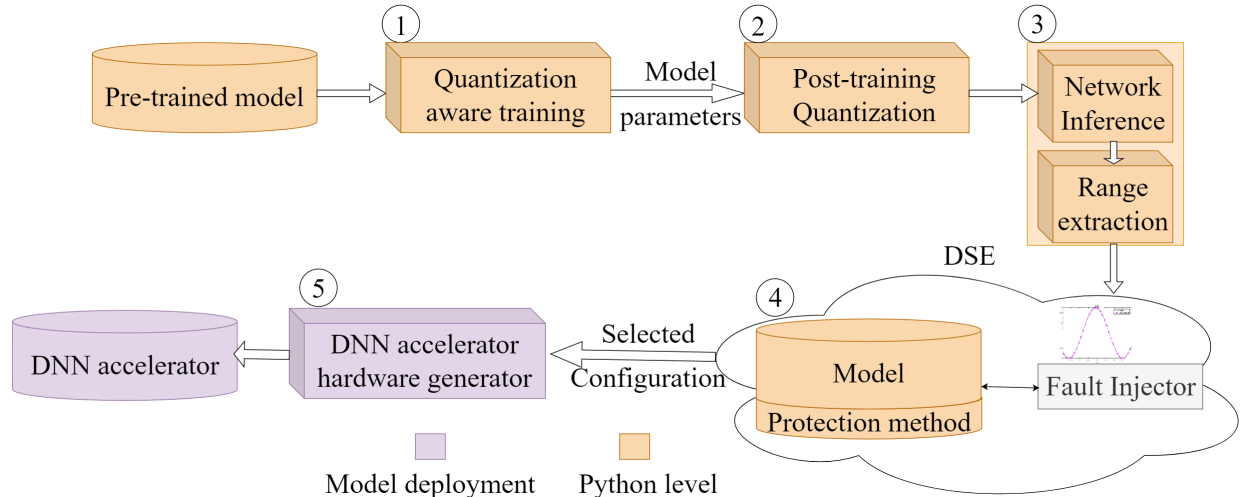
Fig. 2: Proposed methodology flow

these methods can decrease the effect of faults in DNNs, they remove a significant portion of non-zero activations by replacing them with zero, leading to an accuracy drop in high error rates. It is also noteworthy that the mentioned methods do not consider low integer quantization and are mostly working with FP32 and FP16. In this paper, we introduce a novel *lightweight range-checking* circuit that, despite the other works, can consider the maximum values of the layers' activations and replace the out-ranged values with either lower- or upper-bound to avoid fault propagation and also avoid removing a significant portion of non-zero activations by replacing them with only zero. This protection technique is employed in the DNN accelerator hardware generation step of the framework to provide the user with a prototype of the reliable accelerator.

*C. DNN hardware accelerator frameworks*

The advantages of implementing and deploying DNNs on FPGAs are advocated in several recent works. The existing FPGA-based toolchains to map Convolutional Neural Networks (CNNs) are presented in the surveys [16]–[19]. The FINN framework [20] is released by Xilinx for the exploration of quantized CNNs' inference on FPGAs that also provide customized data-flow architectures for each network. Heterogeneous systems are another design strategy in the automated toolchains that propose hardware/software co-design [21]–[23]. In these designs, computational units, e.g., adders and multipliers, are mainly implemented on Programmable Logic (PL) that is controlled by a control unit in a CPU using a dedicated framework, e.g., OpenCL [24]. In this work, we introduce a hardware generation step as part of the framework, to explore DNN inference on an FPGA-based accelerator with a customizable systolic array. It seamlessly integrates with the PYNQ framework [25], leveraging the original PYNQ bootable image. This integration enhances versatility and compatibility, enabling users to implement their network on

different FPGA devices supporting PYNQ. Furthermore, the reconfigurable systolic array implementation introduced in this step provides flexibility and scalability. Users can customize this step to meet their specific network requirements by providing trained parameters and network architectures, resulting in efficient and high-performance DNN inference.

To the best of our knowledge, none of the previous works explored the impact of using different levels of full quantization (weights, activations and biases) of a DNN in the presence of transient faults in the activations on the reliability, accuracy, and delay/resource utilization of the target DNN accelerator.

The approach proposed in this paper goes beyond the state of the art by establishing a fully automated tool for enabling efficient quantization in FPGA-based DNN accelerators aimed at safety-critical applications. The proposed framework contains a high-level simulator to study the impact of quantization on the reliability and accuracy of the network by considering the hardware architecture, with and without protection techniques, followed by an efficient and user-friendly heterogeneous FPGA implementation of the selected DNN configuration.

## III. Proposed Methodology

Fig. 2 illustrates the methodology flow established in the toolchain for reliability and hardware performance analysis of quantized DNN hardware accelerators. This framework takes the DNNs' *Pre-trained model* description as the input. The design, training, and testing of the DNNs are performed in Python. *Quantization-aware training* and the *Post-training quantization*, *Range extraction* and *DSE* steps are seamlessly integrated into the same environment and are responsible for extracting the required data for the hardware generation step. This step is responsible for the hardware implementation of the selected configuration to measure actual hardware parameters like area, latency, etc.

**Step 1: Quantization-aware training.** For this purpose, a full quantization is implemented, targeting all activations, weights, and biases. The framework first takes the description of the network provided by the user and then uses the TFlite library for quantization-aware training. The user can replace their preferred quantization library with the toolchain for this step. The main output of this step is the quantized network's parameters (weights and biases) and network architecture.

**Step 2: Post-training quantization.** In the post-training quantization step, the user can define any further quantization that can be applied to the network with a negligible accuracy loss depending on the level of the quantization. This framework supports quantizing the network down to 4-bit INT. The output accuracy of the generated network is also provided at this step and is kept as a baseline for the further steps of the methodology. For this step, the following algorithm is applied to the network parameters:

The mapping equation is defined as:

$$\tilde{x} = \text{clamp}\left(\left\lfloor \frac{x}{S} \right\rfloor + Z; q_{\min}, q_{\max}\right)$$

$$S = \frac{x_{\max} - x_{\min}}{2^b - 1}$$

Where Z is the offset defined as zero-point, $x_{\max}$ and $x_{\min}$ represent the maximum and the minimum value in the vector. The quantization range $[q_{\min}, q_{\max}]$ is determined by the bit-width. We focus solely on uniform unsigned symmetric quantization, as it is the most commonly employed quantization setup. Hence, $q_{\min}$ is equal to 0, and $q_{\max}$ is equal to $2^b - 1$, where $b$ denotes the bit-width, determining the number of integer grids.

**Step 3: Inference and range extraction.** In this step, after running the inference, the ranges of the activations are extracted for evaluation and reliability study. The ranges are extracted based on the set of validation data, and then the framework extracts the next set of ranges for each layer based on the test data and validates the extracted data correspondingly.

**Step 4: Design Space Exploration.**

**Step 4-A: Fault simulation.** Reliability analysis relies on a Fault Injection (FI) in a *systolic-array-based simulation* of the network in Python, assuming the single bit-flip faults in the activations. While the multiple-bit fault model is more accurate, it requires a prohibitively large number of fault combinations to be considered. Fortunately, it has been shown that high fault coverage obtained using the single-bit model results in a high fault coverage of multiple-bit faults [26]. Therefore, a vast majority of practical FI and test methods are based on the single-bit fault assumption. However, this framework is capable of applying multiple-bit-flips as a fault model depending on the user demand.

The reliability analysis step applies the accuracy drop comparison of the network-under-test as one of the assessment metrics. In addition, the framework assesses the reliability of the DNN by comparing the output probability vector of the golden run (i.e. the DNN that behaves as expected, without

faults) and the faulty run (i.e. the DNN that includes the fault). These metrics involve the SDC (Silent Data Corruption) rate. Specifically, one of the two metrics is "absolute", and the other one is "relative". The SDC rate is defined as the proportion of faults that caused misclassification in comparison with the golden model.

- **SDC-1:** Fault caused a misclassification in the top-ranked output class.
- **SDC-5:** Fault caused the top-ranked element not to exist in the top-5 predicted output classes.
- **SDC-10%:** Fault caused a variation in the output confidence score of the top-ranked output class more than 10% compared to the golden model.

After choosing the preferred quantization in **Step 2**, the designer can go through the systolic-array-based fault injector provided for the reliability evaluation of the Quantized DNN (QDNN). The final design is fed to the next step *hardware generator* for the DNN hardware accelerator generation and hardware performance evaluation process.
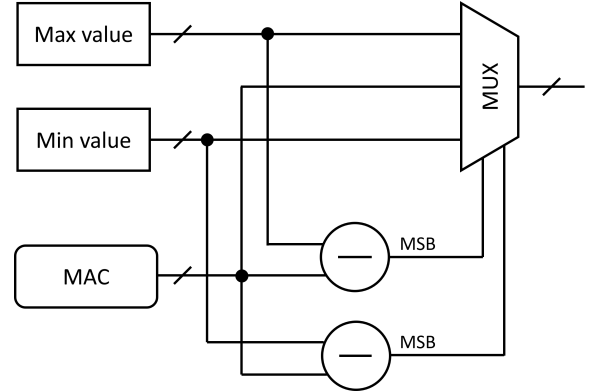


Fig. 3: Proposed lightweight mitigation technique

**Step 4-B: Fault mitigation.** Analyzing the output values of the network's intermediary layers post-training reveals identifiable upper and lower bounds for the neuron's output values. Leveraging this characteristic, we can ensure that any out-of-range outputs are reassigned to the respective upper or lower-bound values. This approach can be effectively implemented using specialized hardware units, as outlined below.

*Out-range Error Detection*: If the neuron's output value exceeds the predetermined upper or lower bound, it indicates a fault in the neuron's input values. To address this, a comparison is made between the neuron's output value and the two pre-established threshold values. For effective error detection, this paper introduces the following strategy.

For each layer, we store two values of upper bound and lower bound as the reference threshold for the out-ranged values. The output of the MAC (Multiply-Accumulate) unit is compared with the threshold values using two subtractors (negative values indicate that the output is beyond the threshold). The result of this comparison defines the final output (Fig. 3). The general overhead of this mitigation technique is two stored values for each layer, and two subtractors to compare
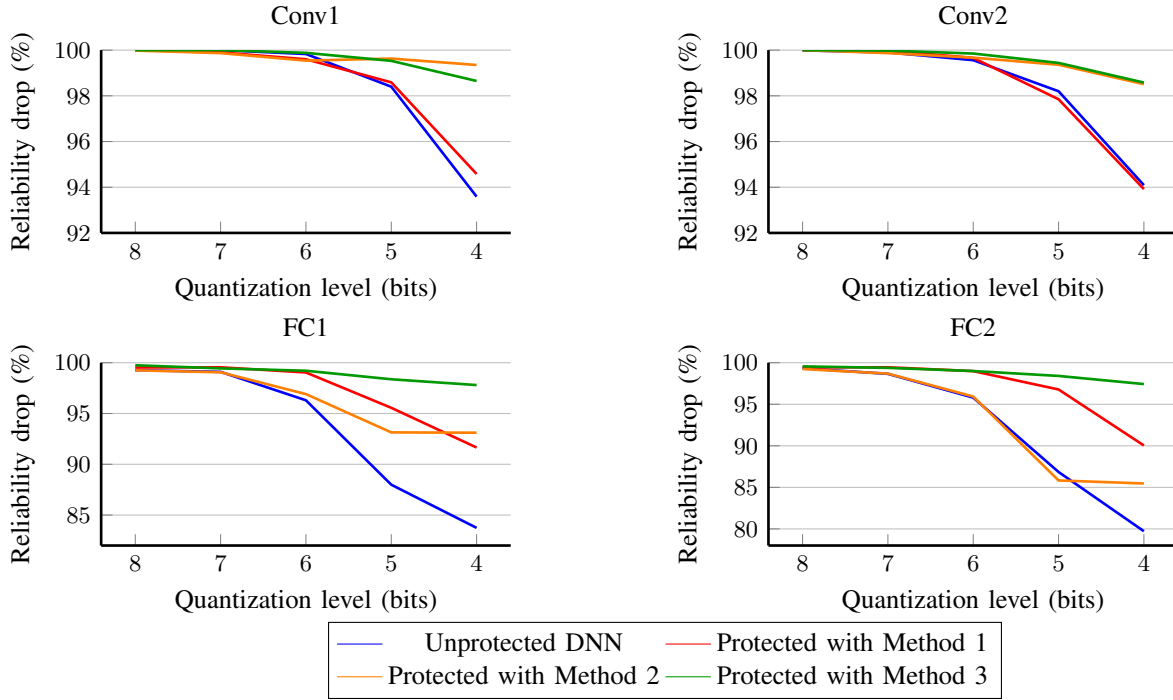
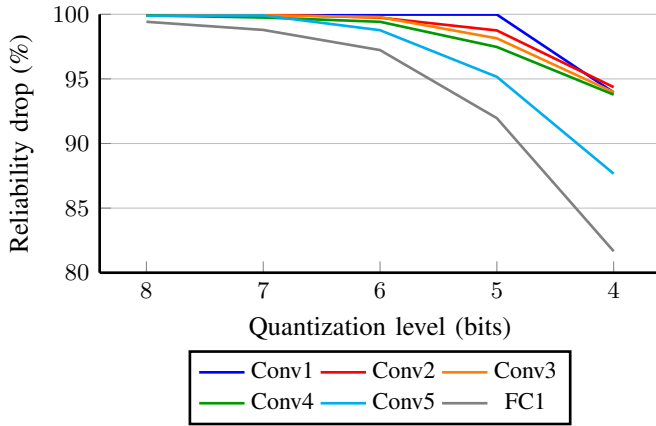Fig. 4: Lenet-5 layer-level reports of reliability drop (based on FI for different quantized networks)



Fig. 5: AlexNet layer-level reports of reliability drop (%) based on different quantization levels (unprotected design)

the MAC output value with the range threshold values and provide the select signal for the MUX to make the decision.

Three variations of this protection technique were implemented in the software to provide users with insights into the reliability enhancements this framework offers:

1) **Method 1:** When out-of-range value is detected it is replaced by the lower bound (min value).
2) **Method 2:** When out-of-range value is detected it is replaced by the upper bound (max value).
3) **Method 3:** When out-of-range value is detected it is replaced by either lower or upper bound depending on

the sign of the MAC output.

This protection technique is designed for easy replacement with any other protection methods (i.e. FT-ClipAct [15]) within this framework toolchain without compromising the overall versatility of the framework.

**Step 5: Hardware generation.**

At this step, a systolic-array-based QDNN accelerator for FPGA SoC is generated based on the parameters of the quantized network provided by **Step 4** to assess hardware utilization and requirements.

The following tasks are executed at this step:

1. Network parameters are analyzed to determine the size of the systolic array, bit precision, and AXI bus bandwidth for data transfer. This analysis takes into account the number of kernels and feature map sizes. The goal is to optimize hardware accelerator performance for the generated network and improve overall efficiency.

2. The board is configured with the PYNQ bootable image. PYNQ provides Python and Jupyter Notebook support to AMD-Xilinx embedded devices. Included Python APIs allow to control both processing system and programmable logic (FPGA). PYNQ setup was selected to provide the users with a familiar interactive Python environment.

3. Network weights and biases are loaded on the board as NumPy array files. The network is described using a provided Python package that interfaces with the accelerator.

4. FPGA is configured from the Jupyter Notebook with the generated accelerator. Then, inference can be run using the provided input data.
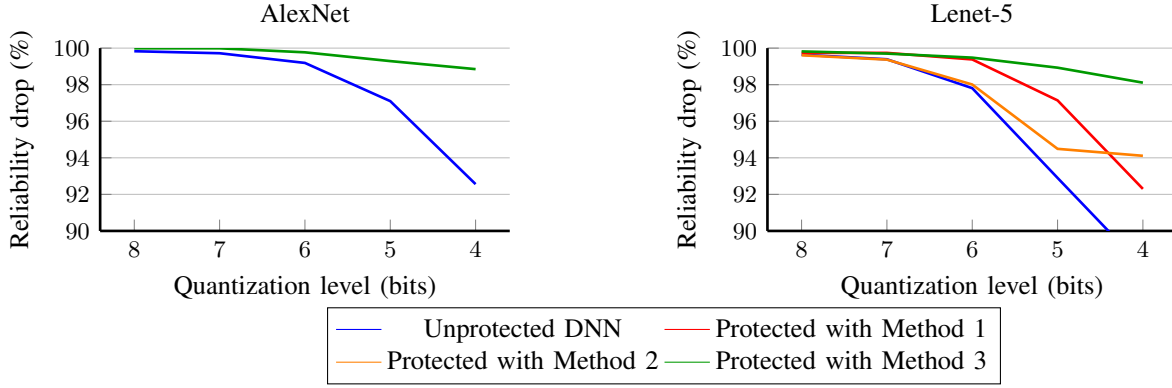
Fig. 6: Model-level reports of reliability drop (%) based on different quantization degrees for AlexNet (left) and LeNet-5 (right)

TABLE I: Lenet-5 layer-level reports of fault criticality (%) based on FI for different quantized networks

| % of critical faults | Unprotected | | | | | Protected with Method 1 | | | | | Protected with Method 2 | | | | | Protected with Method 3 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Lenet-5 | 8 bit | 7 bit | 6 bit | 5 bit | 4 bit | 8 bit | 7 bit | 6 bit | 5 bit | 4 bit | 8 bit | 7 bit | 6 bit | 5 bit | 4 bit | 8 bit | 7 bit | 6 bit | 5 bit | 4 bit |
| conv1 | 0.31 | 0.52 | 1.37 | 3.27 | 9.12 | 0.01 | 0 | 0.49 | 2.82 | 9.06 | 0.3 | 0.6 | 1.45 | 1.69 | 3.76 | 0 | 0 | 0.37 | 1.46 | 3.49 |
| conv2 | 0.29 | 0.46 | 1.33 | 3.62 | 9.38 | 0.07 | 0.08 | 0.84 | 3.42 | 8.49 | 0.21 | 0.57 | 1.27 | 2.45 | 4.71 | 0.07 | 0.08 | 0.51 | 1.71 | 4.08 |
| fc1 | 1.67 | 2.03 | 5.65 | 14.88 | 21.15 | 1.04 | 0.9 | 2.14 | 6.67 | 11.21 | 1.72 | 1.78 | 4.91 | 9.23 | 11.13 | 0.82 | 1.18 | 1.82 | 3.2 | 4.53 |
| fc2 | 1.6 | 2.41 | 5.88 | 16.31 | 25.5 | 1.24 | 1.23 | 1.98 | 4.79 | 13.68 | 1.59 | 2.22 | 5.94 | 17.42 | 19.41 | 0.97 | 1.26 | 2.24 | 3.09 | 5.07 |

TABLE II: AlexNet layer-level reports of fault criticality (%) based on FI for different quantized networks

| % of critical faults | Unprotected | | | | | Protected with Method 3 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| AlexNet | 8 bit | 7 bit | 6 bit | 5 bit | 4 bit | 8 bit | 7 bit | 6 bit | 5 bit | 4 bit |
| conv1 | 0.5 | 0.79 | 1.76 | 4.03 | 8.81 | 0.05 | 0.06 | 0.52 | 1.87 | 3.56 |
| conv2 | 0.58 | 1.05 | 1.35 | 1.66 | 4.11 | 0.03 | 0.03 | 1.035 | 1.39 | 3.31 |
| conv3 | 1.46 | 1.47 | 5.11 | 11.48 | 23.91 | 0.07 | 0.08 | 1.14 | 1.29 | 4.38 |
| conv4 | 0.99 | 1.63 | 2.46 | 7.13 | 14.26 | 0.03 | 0.04 | 1.30 | 4.13 | 5.17 |
| conv5 | 0.90 | 2.10 | 3.69 | 7.82 | 14.31 | 0.04 | 0.09 | 1.61 | 3.44 | 5.17 |
| fc1 | 3.02 | 4.95 | 8.15 | 16.38 | 31.19 | 0.14 | 0.20 | 1.90 | 5.11 | 8.66 |

TABLE III: SDC report for two unprotected Lenet-5 examples with different quantization levels

| Metric (%) | 16-bit | 8-bit |
|---|---|---|
| SDC-1 | 3.18 | 5.24 |
| SDC-5 | 28.04 | 37.26 |
| SDC-10% | 14.30 | 17.65 |

## IV. EXPERIMENTAL RESULTS

### A. Experimental setup

Two networks are studied in this work: Lenet-5 and AlexNet. Lenet-5 is trained on the MNIST dataset, and AlexNet is trained on the CIFAR-10 dataset. Both networks are trained according to the **Step 1** methodology using quantization-aware training. Lenet-5 is trained using 16-bit INT data type, AlexNet is trained using 8-bit INT. For the study, different levels of quantization are applied in the **Step 2** using post-training quantization.

Simulations are performed on 2 × Intel Xeon Gold 6148 2.40 GHz (40 cores, 80 threads per node) with 96 GB RAM. To speed up the simulation process, the framework supports multi-thread parallelism.

To show the hardware characteristics of the output QDNN, studied networks are implemented on the Zynq UltraScale+ ZCU104 Evaluation Board (xczu7ev-ffvc1156-2-e).

### B. Fault simulator

The fault simulator that is used in **Step 4** calculates the sufficient number of faults required for the reliability analysis. QDNNs generated by **Step 2** are validated by means of fault injection over the test set.

*Random fault injection.* According to the adopted fault model, a random single bit-flip is injected into a random activation in a random layer of the network, and the whole test set is fed to the network to obtain the accuracy of the network. This process is repeated several times to reach an acceptable confidence level, which depends on the number of neurons and data representation bit length based on [27]. This work provides an equation to reach 95% confidence level and 1% error margin. The framework adopts the formula presented in this work and provides a sufficient number of repetitions required for reliability analysis.

### C. Validation results

The accuracy results for the quantized networks are reported in Table IV. Further, fault injection is applied on each network automatically as part of the defined configuration of the framework, and reliability drop and fault criticality are

TABLE IV: Model-level design space exploration results for Lenet-5 and AlexNet

| Network | BP | GIOPS | Resource utilization | | | Accuracy, % | Reliability improvement, % | HW utilization (LUT) | | | Fault criticality improvement, % | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | LUT | FF | DSP | | | M1 | M2 | M3 | M1 | M2 | M3 |
| Lenet-5 | 16 | 0.058 | 5298 | 12,892 | 9 | 95.41 | — | 144 | 144 | 576 | — | — | — |
| | 8 | 0.079 | 3475 | 7003 | 9 | 94.02 | 64.33 | 72 | 72 | 288 | 57.78 | 8.75 | 65.61 |
| | 7 | — | — | — | — | 93.93 | 67.95 | 68 | 68 | 135 | 71.24 | 14.95 | 67.74 |
| | 6 | — | — | — | — | 93.52 | 71.90 | 63 | 63 | 99 | 57.25 | 6.16 | 65.91 |
| | 5 | — | — | — | — | 92.49 | 81.17 | 68 | 68 | 81 | 36.30 | 31.34 | 66.86 |
| | 4 | 0.087 | 2114 | 3865 | 9 | 89.65 | 81.17 | 36 | 36 | 63 | 25.85 | 44.92 | 69.20 |
| AlexNet | 16 | 0.338 | 16,654 | 35,503 | 64 | — | — | 1024 | 1024 | 2048 | — | — | — |
| | 8 | 0.465 | 12,138 | 20,539 | 64 | 73.03 | 92.96 | 512 | 512 | 1024 | — | — | 94.27 |
| | 7 | — | — | — | — | 72.26 | 89.79 | 480 | 480 | 960 | — | — | 95.19 |
| | 6 | — | — | — | — | 72.11 | 73.72 | 448 | 448 | 704 | — | — | 58.59 |
| | 5 | — | — | — | — | 70.69 | 66.32 | 480 | 480 | 576 | — | — | 54.13 |
| | 4 | 0.562 | 6428 | 10,067 | 64 | 69.15 | 78.07 | 256 | 256 | 448 | — | — | 60.08 |

reported in Fig. 4 and Table I for the Lenet-5 and in Fig. 5 and Table II for AlexNet. *Reliability drop* is defined as the percentage of accuracy loss in the presence of the faults in the activations in a systolic-array-based simulation model of the network. *Fault criticality* is defined as percentages of the faults that show a negative impact on the network accuracy and lead to misclassification. In Fig. 4 and Table I, the results for all versions of the proposed protection technique are documented for Lenet-5. Table II, only the network protected with Method 3 is compared with the unprotected network for AlexNet, and in Fig. 5 the reports the areliability drop without the protection techniques to show the impact of faults in activations, on different quantization level and layers of an AlexNet network. primarily due to space limitations within the paper.

From the previous works [11], it is evident that the reduction in memory size and quantization can lead to enhanced resilience and mitigate the impact of weight faults due to a reduced memory footprint. However, according to the presented charts, quantization may simultaneously heighten the network's vulnerability to faults in activations and logic. This is particularly crucial in lower precision networks, where even minor bit alterations can have significant ramifications. That is why reliability studies in the DNNs should be done for each QDNN to ensure the impact of quantization on the network's reliability.

Fig. 4 shows that protection Method 3 is capable of improving the reliability of the network in the presence of a fault for more than 34.23% in the worst case for Lenet-5. These numbers are calculated based on the following equation:

$$\% \text{ of Improvement} = \left( \frac{\text{New Value} - \text{Old Value}}{\text{Old Value}} \right) \times 100$$

The same results are reported for AlexNet in Table IV, which shows an improvement of more than 51.79% in the worst case. Improvements in fault criticality for both networks at the model level are also reported in Table IV, which demonstrates the positive impact of the protection technique on reducing the criticality of faults in both networks. These data also showcase the increasing fault criticality in different networks by increasing the level of quantization. Based on the results reported in Table IV, protection Method 3, which shows the

best results for improving reliability among all of the proposed protection techniques, introduces less than 10% overhead compared to the LUTs required for the unprotected network implementation. Meanwhile, full protection of the network with TMR (Triple Module Redundancy) introduces more than 200% hardware overhead.

The fault injection procedure is performed for different quantizations and different versions of the proposed protection technique, and the accuracy drop, due to quantization and fault injection, is profiled. Further, in Table III, SDC metrics of two examples of quantized Lenet-5 are reported. It can be seen that these two networks are susceptible to injected faults. Specifically, the SDC-10% and SDC-5 are very high: on average, about 3.18% of the time the faulty inference misclassified the input in the 16-bit network and 5.24% in the 8-bit network; furthermore, in 28.04% cases for the 16-bit network and 37.26% cases for the 8-bit network, the expected class is not even in the TOP-5 predictions. In addition, it can be observed that the 16-bit quantized network shows better performance in the presence of faults compared to the 8-bit network. In general, these results show that the DNNs used in this experiment are not suitable for a safety-critical application.

Hardware resource utilization and inference latency in GIOPS (Giga Integer Operations Per Second) for different quantization levels are reported in Table IV alongside accuracy, reliability improvement due to the quantization, and hardware overhead and fault criticality improvement for fault mitigation techniques. These results of model-level design space exploration are provided for the user to understand the trade-off between reliability, accuracy, and required computational resources.

## V. CONCLUSION

This paper presents a comprehensive methodology for exploring and enabling a holistic assessment of the trilateral impact of quantization on model accuracy, activation fault reliability, and hardware efficiency. A fully automated framework is introduced that is capable of applying various quantization techniques, fault injection, and hardware implementation, thus enabling the measurement of crucial hardware parameters like area and latency. Moreover, this paper proposes a novel lightweight protection technique integrated within the

framework to ensure the dependable deployment of the final systolic-array-based FPGA implementation. The experiments on established benchmarks demonstrate the analysis flow and the profound implications of quantization on reliability, hardware performance, and network accuracy, particularly concerning the transient faults in the network's activations.

## VI. Acknowledgement

## References

[1] M. Taheri, "Dnn hardware reliability assessment and enhancement," *27th IEEE European Test Symposium (ETS).*, May 2022.

[2] A. Gholami, S. Kim, Z. Dong, Z. Yao, M. W. Mahoney, and K. Keutzer, "A survey of quantization methods for efficient neural network inference," *arXiv preprint arXiv:2103.13630*, 2021.

[3] M. Riazati, M. Daneshtalab, M. Sjödin, and B. Lisper, "Autodeephls: Deep neural network high-level synthesis using fixed-point precision," in *2022 IEEE 4th International Conference on Artificial Intelligence Circuits and Systems (AICAS)*. IEEE, 2022, pp. 122–125.

[4] M. Taheri, M. Riazati, M. H. Ahmadilivani, M. Jenihhin, M. Daneshtalab, J. Raik, M. Sjödin, and B. Lisper, "Deepaxe: A framework for exploration of approximation and reliability trade-offs in dnn accelerators," in *2023 24th International Symposium on Quality Electronic Design (ISQED)*. IEEE, 2023, pp. 1–8.

[5] I. Choi, J.-Y. Hong, J. Jeon, and J.-S. Yang, "Rq-dnn: Reliable quantization for fault-tolerant deep neural networks," in *2023 60th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2023, pp. 1–2.

[6] M. H. Ahmadilivani, M. Taheri, J. Raik, M. Daneshtalab, and M. Jenihhin, "A systematic literature review on hardware reliability assessment methods for deep neural networks," *arXiv preprint arXiv:2305.05750*, 2023.

[7] M. Taheri, M. Taheri, and A. Hadjahmadi, "Noise-tolerance gpu-based age estimation using resnet-50," *arXiv preprint arXiv:2305.00848*, 2023.

[8] A. Siddique, K. Basu, and K. A. Hoque, "Exploring fault-energy trade-offs in approximate dnn hardware accelerators," in *2021 22nd International Symposium on Quality Electronic Design (ISQED)*. IEEE, 2021, pp. 343–348.

[9] 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.

[10] P. M. Basso, F. F. dos Santos, and P. Rech, "Impact of tensor cores and mixed precision on the reliability of matrix multiplication in gpus," *IEEE Transactions on Nuclear Science*, vol. 67, no. 7, pp. 1560–1565, 2020.

[11] F. Libano, P. Rech, B. Neuman, J. Leavitt, M. Wirthlin, and J. Brunhaver, "How reduced data precision and degree of parallelism impact the reliability of convolutional neural networks on fpgas," *IEEE Transactions on Nuclear Science*, vol. 68, no. 5, pp. 865–872, 2021.

[12] F. Libano, B. Wilson, M. Wirthlin, P. Rech, and J. Brunhaver, "Understanding the impact of quantization, accuracy, and radiation on the reliability of convolutional neural networks on fpgas," *IEEE Transactions on Nuclear Science*, vol. 67, no. 7, pp. 1478–1484, 2020.

[13] R. T. Syed, M. Ulbricht, K. Piotrowski, and M. Krstic, "Fault resilience analysis of quantized deep neural networks," in *2021 IEEE 32nd International Conference on Microelectronics (MIEL)*. IEEE, 2021, pp. 275–279.

[14] A. P. Arechiga and A. J. Michaels, "The effect of weight errors on neural networks," in *2018 IEEE 8th Annual Computing and Communication Workshop and Conference (CCWC)*. IEEE, 2018, pp. 190–196.

[15] L.-H. Hoang, M. A. Hanif, and M. Shafique, "Ft-clipact: Resilience analysis of deep neural networks and improving their fault tolerance using clipped activation," in *2020 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2020, pp. 1241–1246.

[16] S. I. Venieris, A. Kouris, and C.-S. Bouganis, "Toolflows for mapping convolutional neural networks on fpgas: A survey and future directions," *arXiv preprint arXiv:1803.05900*, 2018.

[17] K. Guo, S. Zeng, J. Yu, Y. Wang, and H. Yang, "[dl] a survey of fpga-based neural network inference accelerators," *ACM Transactions on Reconfigurable Technology and Systems (TRETS)*, vol. 12, no. 1, pp. 1–26, 2019.

[18] K. Abdelouahab, M. Pelcat, J. Serot, and F. Berry, "Accelerating cnn inference on fpgas: A survey," *arXiv preprint arXiv:1806.01683*, 2018.

[19] R. S. Molina, V. Gil-Costa, M. L. Crespo, and G. Ramponi, "High-level synthesis hardware design for fpga-based accelerators: Models, methodologies, and frameworks," *IEEE Access*, vol. 10, pp. 90 429–90 455, 2022.

[20] Y. Umuroglu, N. J. Fraser, G. Gambardella, M. Blott, P. Leong, M. Jahre, and K. Vissers, "Finn: A framework for fast, scalable binarized neural network inference," in *Proceedings of the 2017 ACM/SIGDA international symposium on field-programmable gate arrays*, 2017, pp. 65–74.

[21] S. I. Venieris and C.-S. Bouganis, "fpgaconvnet: Mapping regular and irregular convolutional neural networks on fpgas," *IEEE transactions on neural networks and learning systems*, vol. 30, no. 2, pp. 326–342, 2018.

[22] A. Ghaffari and Y. Savaria, "Cnn2gate: Toward designing a general framework for implementation of convolutional neural networks on fpga," *arXiv preprint arXiv:2004.04641*, 2020.

[23] P. G. Mousouliotis and L. P. Petrou, "Cnn-grinder: from algorithmic to high-level synthesis descriptions of cnns for low-end-low-cost fpga socs," *Microprocessors and Microsystems*, vol. 73, p. 102990, 2020.

[24] J. E. Stone, D. Gohara, and G. Shi, "Opencl: A parallel programming standard for heterogeneous computing systems," *Computing in science & engineering*, vol. 12, no. 3, p. 66, 2010.

[25] E. Wang, J. J. Davis, and P. Y. Cheung, "A pynq-based framework for rapid cnn prototyping," in *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 2018, pp. 223–223.

[26] M. Bushnell and V. Agrawal, *Essentials of electronic testing for digital, memory and mixed-signal VLSI circuits*. Springer Science & Business Media, 2004, vol. 17.

[27] R. Leveugle, A. Calvez, P. Maistri, and P. Vanhauwaert, "Statistical fault injection: Quantified error and confidence," in *2009 Design, Automation & Test in Europe Conference & Exhibition*. IEEE, 2009, pp. 502–506.