

Article Information

Title	Approximated Triple Modular Redundancy of Convolutional Neural Networks Based on Residual Quantization
Authors	Yamato Saikawa, Yoichi Tomioka
Citation	2024 IEEE 17th International Symposium on Embedded Multicore/Many-core Systems-on-Chip (MCSoc), Kuala Lumpur, Malaysia, 2024, pp. 302-309, doi: 10.1109/MCSoc64144.2024.00057.
Copyright	© 2024 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.
Note	<p>This is the author's accepted version of the paper published in <i>Proceedings of IEEE International Symposium on Embedded Multicore/Many-core Systems-on-Chip (2024)</i>.</p> <p>The final version is available at IEEE Xplore DOI: https://doi.org/10.1109/MCSoc64144.2024.00057</p>

Approximated Triple Modular Redundancy of Convolutional Neural Networks Based on Residual Quantization

1st Yamato Saikawa and 2nd Yoichi Tomioka
Graduate School of Computer Science and Engineering
The University of Aizu
Aizu-Wakamatsu City, Fukushima, Japan
{d8261107,ytomioka}@u-aizu.ac.jp

Abstract—Triple modular redundancy (TMR) is a fault-tolerant system that detects faults and allows processing to continue by employing three identical modules. However, TMR introduces challenges such as tripled circuit area and increased power consumption. In this paper, we propose an approximated TMR (ATMR), which consists of the original module and two smaller modules achieved through quantization, specifically tailored for convolutional neural networks (CNNs). For evaluation, we conduct experiments using ResNet-20 to simulate fault recovery through linear and residual quantization when the original module fails. The results demonstrate that ATMR maintains high inference accuracy across many layers, even when employing low-bit linear or residual quantization for recovery.

Index Terms—convolutional neural network, quantization, triple modular redundancy, fault detection, restoration, fault tolerance

I. INTRODUCTION

Edge AI systems based on deep neural networks (DNNs) have started being used in various fields and affect our lives closely. The DNNs that are most utilized as current edge AIs are based on Artificial neural networks (ANNs). It mainly calculates matrix multiplication with real or integer values and achieves high inference accuracy for data analysis, classification, pattern recognition, etc.

Convolutional neural networks (CNNs) [1] can classify multidimensional data more accurately than general ANNs. It is reported that the CNNs of deeper layers have more problems with gradient loss or explosion. Therefore, finding the appropriate parameters for such CNNs is difficult through back-propagation-based training even though they have the potential to improve accuracy. To solve the problems, Residual Network (ResNet) [2] is conceived. This method features a skip connection that has connected identity mapping and residual mapping, and it is called residual learning. Thanks to it, CNNs are able to consist of many layers and obtain very high inference accuracy. However, CNN inference is computationally expensive and requires a dedicated circuit for real-time processing. Such a circuit requires high power consumption and a larger circuit area to achieve higher performance since to

achieve faster CNN inference, it is necessary to achieve highly parallel computation. In particular, dot products, including convolution operations, request many resources such as logic gates, registers, etc. Therefore, CNN accelerators of low power and small size are needed in order for CNNs to be utilized on edge devices that have serious restrictions such as power consumption and heat dissipation.

To solve the size and power consumption problems of CNNs, there have been works of model compression and several of them have shown high compression ratios and inference accuracy. For example, the method called quantization [3]–[5] can reduce the original bit width to lower. Keeping the structure of the original model, the size and computing resources are able to be compressed owing to that.

On the other hand, there is report [6] that circuits may have crises, such as unintentional behavior due to aging or physical failure. Especially, malfunction of automatic driving systems, robots, medical devices, and so on have possible serious accidents affect human life. Thus, AI in safety-critical edge systems should be realized to detect malfunctions and prevent accidents in advance.

As the existing method to keep high reliability, there is triple modular redundancy (TMR) that consists of three modules of the same circuit. It can detect the fault and continue working normally when one of the three is out of order as it gathers outputs from three circuits and selects their majority. Although TMR can get robustness of systems, it is difficult to use large-scale systems in edge devices like chips with small sizes because implementation costs such as power consumption and circuit area increase by a factor of three or more. In the current situation, the TMR system in a CNN is hard to adapt to the implementation of small-size chips. Thus, we aim to relax the problems and achieve fault tolerance with a smaller area and lower power.

From the above, we propose an approximated TMR (ATMR) based on quantization for CNNs. In this method, two redundant modules are replaced in quantized them, and they play a role in fault detection and substitute of inference if the original module is malfunctioning. Moreover, we also propose restoration by two methods that are based on linear

quantization and residual quantization. Hence, it is considered that the problems that the conventional TMR has can be relaxed while the proposed method keeps high reliability.

Our contributions to this thesis are summarized as follows:

- We propose ATMR(8, m , m) that consists of the original convolutional layer and two m -bit quantized convolutional layers. This work extends approximated dual modular redundancy (ADMR) [7], which is based on linear quantization. While ADMR cannot restore outputs upon failure, the proposed ATMR enables continuous inference even when faults are detected. By partially applying residual quantization, ATMR achieves an accuracy comparable to a non-faulty model, despite requiring less computational cost than the conventional TMR approach.
- We perform output restoration experiments of ATMR(8, m , m) which consists of 8-bit and m -bit convolutional layers ($m = \{3, 4, 5, 6, 7\}$) for 8-bit ResNet-20 trained with CIFAR-10.
- In software simulation, the output restoration of ResNet-20 based on linear quantization achieves that ATMR(8,3,3) and ATMR(8,4,4) recover the inference accuracy to near 91.8% which the original has when their area under the curve (AUC) of fault detection is 1 in many layers. Moreover, we demonstrate that residual quantization helps improve the accuracy in layers where the accuracy drops by linear quantization-based ATMR.

II. RELATED WORKS

Reference [8] proposes a fault-tolerance method based on N-version programming. In this approach, multiple neural networks are independently trained with different parameters, architectures, or subsets of the training data and executed in parallel. The study evaluated the method using the MNIST and CIFAR-10 datasets, showing that parallel networks achieved higher accuracy than individual networks. However, the method faces challenges, such as performance degradation caused by permanent circuit faults. Reference [9] presents a technique for generating quantized neural networks with low-bit precision and performing multiply-accumulate operations with a guard bit to enhance fault tolerance. Similarly, Reference [10] proposes constructing dual modular redundancy (DMR) at the neuron level to mitigate performance degradation under single-bias attacks. However, these methods have limitations: they cannot address faults affecting entire layers because they only duplicate individual neurons, nor can they handle faults other than increased delays in multiply-accumulate operations. Reference [11] introduces an approximated DMR unit that employs ternary convolutional layers and random forests to significantly reduce circuit area and computational costs while detecting a variety of faults. However, this method is not directly applicable to general CNN models with wide bit widths, limiting its versatility in broader applications.

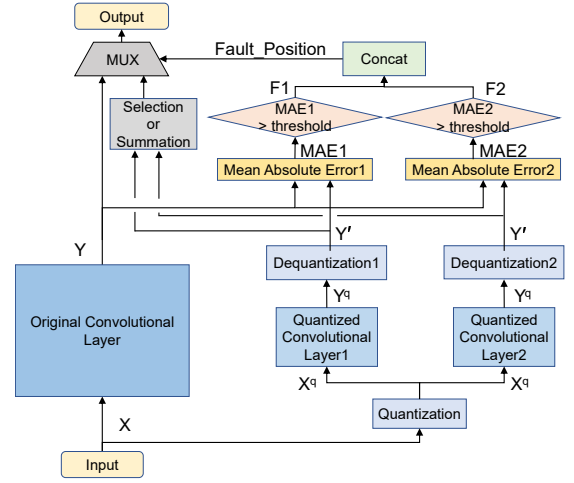


Fig. 1. Proposed ATMR.

III. PROPOSED ATMR METHOD

A. ATMR Outline

The outline of the proposed ATMR is shown in Figure 1. \mathbf{X} is a input tensor that consists of $C_{in} \times H \times W$, where C_{in} is the number of input channel, H is the height of each channel and W is the width.

First of all, an input \mathbf{X} enters the Original Convolutional Layer and the Quantization module. The Original Convolutional Layer is the main layer to infer and perform convolution operations with original-bit input and weight. After the calculation, it outputs \mathbf{Y} that consists of $C_{out} \times H \times W$ elements. It is denoted by

$$\mathbf{Y} = \text{Conv}(\mathbf{X}, \mathbf{W}). \quad (1)$$

Defining c , h , and w as the indices of tensor elements, respectively, the output of $\mathbf{Y}_{c,w,h}$ is calculated as follows:

$$\begin{aligned} \mathbf{Y}_{c,h,w} &= \sum_{i=0}^{C_{in}-1} \sum_{j=0}^{K-1} \sum_{k=0}^{K-1} \mathbf{X}_{i,h+j-\lfloor \frac{K}{2} \rfloor, w+k-\lfloor \frac{K}{2} \rfloor} \cdot \mathbf{W}_{c,i,j,k}, \quad (2) \end{aligned}$$

where K is the size of the kernel which the layer has as weights, \mathbf{W} is a tensor as a weight that consists of $C_{out} \times C_{in} \times K \times K$ elements. \mathbf{W}_c represents c -th weight filter in the weight tensor \mathbf{W} .

The quantization module is a module to quantize \mathbf{X} with the original bits to \mathbf{X}^q with the lower bits. The module has α as the scaling factor and γ as the zero-point to compress to a specified bit width with a minimum value of 0. Moreover, the quantization is performed for each channel. Thus, \mathbf{X}^q , α and γ are calculated as follows:

$$\mathbf{X}^q = \text{round}(\alpha(\mathbf{X} - \gamma)), \quad (3)$$

$$\alpha = \frac{2^m - 1}{\max\{\mathbf{X}\} - \min\{\mathbf{X}\}} \quad (4)$$

$$\gamma = \min\{\mathbf{X}\}, \quad (5)$$

where m is the bits number after the quantization, \max is the function to return the maximum value from \mathbf{X} and \min is the function to return the minimum value.

After the quantization, \mathbf{X}^q enters two Quantized Convolutional Layers. They are modules to calculate convolution with lower-bit quantized weights and play a role in fault detection and restoration when the inference module is broken. Additionally, they have the same \mathbf{W}^q that is generated by linear quantization based on \mathbf{W} . \mathbf{W}^q and scaling factor β for it are calculated as follows:

$$\mathbf{W}^q = \text{round}(\beta\mathbf{W}) \quad (6)$$

$$\beta = \frac{2^{m-1} - 1}{\max\{|\mathbf{W}|\}} \quad (7)$$

Let \mathbf{Y}^q be the result of the convolutional operation of the two quantization tensors, \mathbf{X}^q and \mathbf{W}^q . It is calculated by

$$\mathbf{Y}^q = \text{Conv}(\mathbf{X}^q, \mathbf{W}^q). \quad (8)$$

\mathbf{Y}^q has different bits from the original layer. The value used for fault detection and output restoration is required to be the same bit width as \mathbf{Y} . Therefore, the approximated value, \mathbf{Y}' is generated by dequantizing in each Dequantization module. The dequantization of \mathbf{X}^q and \mathbf{W}^q can be calculated as follows:

$$\mathbf{X} \approx \frac{\mathbf{X}^q}{\alpha} + \gamma \quad (9)$$

$$\mathbf{W} \approx \frac{\mathbf{W}^q}{\beta} \quad (10)$$

Therefore, the convolutional operation of their approximated tensors also can be calculated as follows:

$$\begin{aligned} \mathbf{Y}_{c,h,w} &\approx \\ &\sum_{i=0}^{C_{in}-1} \sum_{j=0}^{K-1} \sum_{k=0}^{K-1} \left(\frac{\mathbf{X}_{i,h+j-\lfloor \frac{K}{2} \rfloor, w+k-\lfloor \frac{K}{2} \rfloor}^q}{\alpha} + \gamma \right) \cdot \frac{\mathbf{W}_{c,i,j,k}^q}{\beta} \\ &= \frac{1}{\alpha\beta} \sum_{i=0}^{C_{in}-1} \sum_{j=0}^{K-1} \sum_{k=0}^{K-1} \mathbf{X}_{i,h+j-\lfloor \frac{K}{2} \rfloor, w+k-\lfloor \frac{K}{2} \rfloor}^q \cdot \mathbf{W}_{c,i,j,k}^q + \\ &\quad \gamma \sum_{i=0}^{C_{in}-1} \sum_{j=0}^{K-1} \sum_{k=0}^{K-1} \frac{\mathbf{W}_{c,i,j,k}^q}{\beta} \\ &= \frac{1}{\alpha\beta} \mathbf{Y}^q_c + \gamma \cdot \text{sum}(\mathbf{W}_c) \end{aligned} \quad (11)$$

From the above, let α' be $\alpha\beta$ as the new scaling factor and γ' be $\gamma \cdot \text{sum}(\mathbf{W}_c)$ as the new zero point, \mathbf{Y}' as the dequantized tensor is calculated as follows:

$$\mathbf{Y}' = \frac{\mathbf{Y}^q}{\alpha'} + \gamma' \quad (12)$$

Secondly, the fault occurrence of each convolutional layer is predicted, and output restoration is performed. Based on the fault detection method in [7], mean absolute errors (MAEs) between \mathbf{Y} and two \mathbf{Y}' are calculated as follows:

$$MAE = \frac{\|\mathbf{Y} - \mathbf{Y}'\|_1}{C \cdot H \cdot W} \quad (13)$$

After that, the fault is predicted by checking whether the values exceed a threshold. The results of Mean Absolute Error 1 and 2 are denoted by $MAE1$ and $MAE2$, respectively. Each value is fed into a corresponding module that outputs 1 if the input exceeds the threshold or 0 otherwise. This process is formulated as a two-class classification problem, where each layer is classified as either being in a fault state (positive class) or a normal state (negative class) based on its MAE . The MAE between the output \mathbf{Y} from the faulty layer and \mathbf{Y}' from the normal layer is treated as a positive sample, regardless of any inferred class changes due to the fault in that layer. Conversely, the MAE from the normal layer's output is considered a negative sample.

The number of correctly classified positive samples is defined as True Positives (TP), while the number of correctly classified negative samples is defined as True Negatives (TN). Similarly, the number of misclassified positive samples is defined as False Negatives (FN), and the number of misclassified negative samples is defined as False Positives (FP). The True Positive Rate (TPR) and False Positive Rate (FPR) are calculated as follows:

$$TPR = \frac{TP}{FN + TP} \quad (14)$$

$$FPR = \frac{FP}{TN + FP} \quad (15)$$

Additionally, a simulation using several image data is conducted, and based on that, thresholds that satisfy the specified FPR are obtained from the MAE distribution in normal and fault layers.

After the conditional branch, two binary results are combined as $Fault_Position$ with 2 bits to identify fault layers, and it is used as the control signal of the multiplexer. For example, $Fault_Position$ becomes $\{1, 1\}$ when the original Convolutional Layer has a fault, and others are normal. When all layers work normally, $\{0, 0\}$ is output. When one of the Quantized Convolutional Layers is out of order, $\{1, 0\}$ or $\{0, 1\}$ are output. Finally, the multiplexer determines which output of layers to send to next based on the control signal. If the approximated output is determined, a selection of one side or an average of two is used. Furthermore, there are two types of approximated output. The first is to use the tensor calculated by Eq. 12. The second is to use the tensor generated by residual quantization and dequantization.

B. Residual Quantization for Output Restoration

Residual quantization [12] iteratively quantizes the quantization error, and the values are used for narrowing the error range during dequantization to the original value. Thus, the method can approximate values more closely to the original compared to linear quantization in exchange for calculation costs of iterative quantization.

Assuming k iterations and input \mathbf{X} , the quantized $\mathbf{X}^q_{(k)}$, $\alpha_{(k)}$ as the scaling factor and $\gamma_{(k)}$ as the zero point are calculated as follows:

$$\mathbf{X}^q_{(k)} = \text{round}(\alpha_{(k)}(\Delta\mathbf{X}^q_{(k-1)} - \gamma_{(k)})) \quad (16)$$

$$\alpha_{(k)} = \frac{2^m - 1}{\max\{\Delta\mathbf{X}^q_{(k-1)}\} - \min\{\Delta\mathbf{X}^q_{(k-1)}\}} \quad (17)$$

$$\gamma_{(k)} = \min_{\Delta\mathbf{X}^q_{(k-1)}}, \quad (18)$$

where $\Delta\mathbf{X}^q_{(k-1)}$ means a quantization error between $\Delta\mathbf{X}^q_{(k-2)}$ and $\mathbf{X}^q_{(k-1)}$, and it is calculated as follows:

$$\Delta\mathbf{X}^q_{(k)} = \begin{cases} \mathbf{X} & \text{if } k = 0, \\ \Delta\mathbf{X}^q_{(k-1)} - \frac{\mathbf{X}^q_{(k)}}{\alpha_{(k)}} - \gamma_{(k)} & \text{if } k \geq 1 \end{cases} \quad (19)$$

Algorithm 1 Calculation of \mathbf{X}^q with N iterations.

Input: \mathbf{X}

Output: \mathbf{X}^q

```

for  $r = 0$  to  $N$  do
  if  $r = 0$  then
     $\Delta\mathbf{X}^q[r] \leftarrow \mathbf{X}$ 
  else
     $\gamma[r] \leftarrow \min(\Delta\mathbf{X}^q[r-1])$ 
     $\alpha[r] \leftarrow 2^m - 1 / \{\max(\Delta\mathbf{X}^q[r-1]) - \min(\Delta\mathbf{X}^q[r-1])\}$ 
     $\mathbf{X}^q[r] \leftarrow \text{round}(\alpha[r] * \Delta\mathbf{X}^q[r-1] - \gamma[r])$ 
     $\Delta\mathbf{X}^q[r] \leftarrow \Delta\mathbf{X}^q[r-1] - \mathbf{X}^q[r] / \alpha[r] - \gamma[r]$ 
  end if
end for
return  $\mathbf{X}^q$ 

```

Collectively, the above equations from 16 to 19 are executed and the items are obtained such as Algorithm 1. Therefore, assuming that iterations is N , the result of input residual quantization becomes $\mathbf{X}^q = \{\mathbf{X}^q_{(1)}, \mathbf{X}^q_{(2)}, \mathbf{X}^q_{(3)}, \dots, \mathbf{X}^q_{(N-1)}, \mathbf{X}^q_{(N)}\}$, and it is equivalent to linear quantization if $N = 1$.

Furthermore, assuming k iterations and weight \mathbf{W} , the quantized $\mathbf{W}^q_{(k)}$ and $\beta_{(k)}$ as the scaling factor are calculated as follows:

$$\mathbf{W}^q_{(k)} = \text{round}(\beta_{(k)}\Delta\mathbf{W}^q_{(k-1)}) \quad (20)$$

$$\beta_{(k)} = \frac{2^m - 1}{\max\{|\Delta\mathbf{W}^q_{(k-1)}|\}}, \quad (21)$$

where $\Delta\mathbf{W}^q_{(k-1)}$ means a quantization error between $\Delta\mathbf{W}^q_{(k-2)}$ and $\mathbf{W}^q_{(k-1)}$, and it is calculated as follows:

$$\Delta\mathbf{W}^q_{(k)} = \begin{cases} \mathbf{W} & \text{if } k = 0, \\ \Delta\mathbf{W}^q_{(k-1)} - \frac{\mathbf{W}^q_{(k)}}{\beta_{(k)}} & \text{if } k \geq 1 \end{cases} \quad (22)$$

Collectively, the above equations from 20 to 22 are executed and the items are obtained such as Algorithm 2. Therefore, assuming that iterations is N' , the result of residual quantization for weight becomes $\mathbf{W}^q =$

Algorithm 2 Calculation of \mathbf{W}^q with N' iterations.

Input: \mathbf{W}

Output: \mathbf{W}^q

```

for  $r' = 0$  to  $N'$  do
  if  $r = 0$  then
     $\Delta\mathbf{W}^q[r'] \leftarrow \mathbf{W}$ 
  else
     $\beta[r'] \leftarrow 2^{m-1} - 1 / \max(\text{abs}(\Delta\mathbf{W}^q[r'-1]))$ 
     $\mathbf{W}^q[r'] \leftarrow \text{round}(\beta[r'] * \Delta\mathbf{W}^q[r'-1])$ 
     $\Delta\mathbf{W}^q[r'] \leftarrow \Delta\mathbf{W}^q[r'-1] - \mathbf{W}^q[r'] / \beta[r']$ 
  end if
end for
return  $\mathbf{W}^q$ 

```

$\{\mathbf{W}^q_{(1)}, \mathbf{W}^q_{(2)}, \mathbf{W}^q_{(3)}, \dots, \mathbf{W}^q_{(N'-1)}, \mathbf{W}^q_{(N')}\}$. Let \mathbf{Y}^q be the result of the convolutional operation of the two quantization tensors, $\mathbf{X}^q_{(r)}$ and $\mathbf{W}^q_{(r')}$, the calculation of $\mathbf{Y}^q_{(r,r')}$ is the following:

$$\mathbf{Y}^q_{(r,r')} = \text{Conv}(\mathbf{X}^q_{(r)}, \mathbf{W}^q_{(r')}) \quad (23)$$

This operation is performed for all combination of quantized tensors that are in \mathbf{X}^q and \mathbf{W}^q such as Algorithm 3. Thus, \mathbf{Y}^q becomes a tensor that consists of $N \times N' \times C_{out} \times H \times W$.

Algorithm 3 Calculation of \mathbf{Y}^q from N iterations \mathbf{X}^q and N' iterations \mathbf{W}^q

Input: $\mathbf{X}^q, \mathbf{W}^q$

Output: \mathbf{Y}^q

```

for  $r = 1$  to  $N$  do
  for  $r' = 1$  to  $N'$  do
     $\mathbf{Y}^q[r][r'] \leftarrow \text{Conv}(\mathbf{X}^q[r], \mathbf{W}^q[r'])$ 
  end for
end for
return  $\mathbf{Y}^q$ 

```

As same as linear quantization, $\mathbf{Y}^q_{(r)}$ has different bit width from \mathbf{Y} . Therefore, dequantization is performed for the output restoration. The dequantization of $\mathbf{X}^q_{(r)}$ and $\mathbf{W}^q_{(r')}$ calculated by residual quantization can be calculated as follows:

$$\mathbf{X} \approx \sum_{r=1}^N \frac{\mathbf{X}_{(r)}}{\alpha_{(r)}} + \sum_{r=1}^N \gamma_{(r)} + \Delta\mathbf{X}_{(N)} \quad (24)$$

$$\mathbf{W} \approx \sum_{r'=1}^{N'} \frac{\mathbf{W}_{(r')}}{\beta_{(r')}} + \Delta\mathbf{W}_{(N')} \quad (25)$$

Based on Eq. 24 and 25, the approximated value of \mathbf{Y} also can be calculated as follows:

$$\begin{aligned}
\mathbf{Y}_{c,h,w} &\approx \\
&\sum_{r=1}^N \sum_{r'=1}^{N'} \sum_{i=0}^{C_{in}-1} \sum_{j=0}^{K-1} \sum_{k=0}^{K-1} \\
&\left(\frac{\mathbf{X}_{(r),i,h+j-\lfloor \frac{K}{2} \rfloor, w+k-\lfloor \frac{K}{2} \rfloor}^q}{\alpha_{(r)}} + \gamma_{(r)} \right) \cdot \frac{\mathbf{W}_{(r'),c,i,j,k}^q}{\beta_{(r')}} \\
&= \sum_{r=1}^N \sum_{r'=1}^{N'} \frac{1}{\alpha_{(r)}\beta_{(r')}} \sum_{i=0}^{C_{in}-1} \sum_{j=0}^{K-1} \sum_{k=0}^{K-1} \\
&\quad \mathbf{X}_{(r),i,h+j-\lfloor \frac{K}{2} \rfloor, w+k-\lfloor \frac{K}{2} \rfloor}^q \cdot \mathbf{W}_{(r'),c,i,j,k}^q + \\
&\quad \sum_{r=1}^N \gamma_{(r)} \sum_{r'=1}^{N'} \sum_{i=0}^{C_{in}-1} \sum_{j=0}^{K-1} \sum_{k=0}^{K-1} \frac{\mathbf{W}_{(r'),c,i,j,k}^q}{\beta_{(r')}} \\
&= \sum_{r=1}^N \sum_{r'=1}^{N'} \frac{1}{\alpha_{(r)}\beta_{(r')}} \mathbf{Y}^q_{(r,r'),c,h,w} + \sum_{r=1}^N \gamma_{(r)} \cdot \text{sum}(\mathbf{W}_c) \quad (26)
\end{aligned}$$

Therefore, let γ' be $\sum_{r=1}^N \gamma_{(r)} \cdot \text{sum}(\mathbf{W}_c)$ as the new zero point, the dequantization for the output generated by convolutional layer based on residual quantization is calculated as follow:

$$\mathbf{Y}' = \sum_{r=1}^N \sum_{r'=1}^{N'} \frac{1}{\alpha_{(r)}\beta_{(r')}} \mathbf{Y}^q_{(r,r')} + \gamma' \quad (27)$$

IV. RESTORATION EXPERIMENT

A. Experiment Overview

In this experiment, we evaluate the effectiveness of output restoration. Fault generation, detection, and restoration simulations are performed on the software. For fault injection, we employ the Act Bit Flip method [7], which flips a single bit in the output feature maps with a specified probability. In restoration simulation, when a fault is detected, an approximated output based on either linear quantization or residual quantization is adopted. The iteration counts for residual quantization, N and N' , are set to 2. We use a pre-trained ResNet-20 model as our evaluation model. Assuming an edge device implementation, all calculations in the original convolutional layers are performed with 8-bit precision. PyTorch 2.0.1 serves as the deep learning framework, and the model is further quantized to 8-bit precision using PyTorch's Quantization API. The layers in the 8-bit ResNet-20 model are converted to the proposed ATMR, with quantized bit widths m ranging from 3 to 7. The ATMR architecture, which consists of an 8-bit convolutional layer followed by two m -bit quantized convolutional layers, is denoted as ATMR(8, m , m).

Restoration is simulated using an 8-bit ResNet-20 model trained on the CIFAR-10 dataset, achieving an inference accuracy of 91.8%. The model's architecture is shown in Figure 2 (a). The first residual block in Layers 2 and 3 includes downsampling for the input, and its architecture is shown in Figure 2 (b). Other blocks are shown in Figure 2 (c). In this model, the area under the curves (AUCs) and thresholds

corresponding to each layer are acquired from 10,000 images after setting FPR to 0.001. We select one convolutional layer from each layer group and the Stem layer, then evaluate the effects of faults and output restoration by injecting Act Bit Flip faults into the selected layers. Table I presents the changes in accuracy when Act Bit Flip faults occur with probabilities ranging from 1% to 20%. Based on the results, we identify the layers most affected by faults as targets for output restoration. Due to space limitations, in this paper, we present the results of output restoration for the selected layers (*stem.conv1*, *layer1.0.conv1*, *layer2.2.conv1*, *layer3.3.conv1*) where significant inference accuracy degradation is observed by fault injection.

B. Results

First, we present AUCs of fault detection of *stem.conv1* in Table II, and the results of output restoration with linear quantization and residual quantization for *stem.conv1* in Table III and Table IV, respectively. Both restorations of ATMR(8,3,3) cannot reach the original accuracy even when the AUC is 1, as *stem.conv1* is sensitive to faults in terms of inference accuracy, and 3-bit residual quantization is insufficient for recovery given the low values of N and N' . With ATMR(8, 4, 4), the AUC reaches nearly 1 at a fault probability of 11%. The accuracy achieved through linear restoration is only about 2% lower than the original. Residual quantization achieves full recovery when the AUC is 1. This indicates that the bit width required for output restoration in ResNet-20 is at least 4 bits. Both restorations in ATMR(8,5,5) are less than a 10% reduction in inference accuracy since the AUC is 1 at 8%. The restoration by residual quantization performs slightly better for inference accuracy.

Second, we present AUCs for *layer1.0.conv1* in Table V, along with the results of output restoration using linear quantization in Table VI and residual quantization in Table VII. For ATMR(8,3,3), restoration with residual quantization fully recovers the output when the AUC is 1 and performs better than linear quantization. From ATMR(8,4,4), both restorations have consistently maintained over 90%, and they are little difference.

Third, we present AUCs for *layer2.2.conv1* in Table VIII, and the results of output restoration with linear quantization in Table IX and residual quantization in Table X. For ATMR(8,3,3), the AUCs reach 1 at fault probabilities where accuracy degradation becomes significant. Moreover, their TMRs show sufficient inference accuracy for both restorations.

Finally, we present AUCs for *layer3.1.conv1* in Table XI, and the results of output restoration with linear quantization in Table XII and residual quantization in Table XIII. For ATMR(8,3,3), linear quantization provides better restoration accuracy than residual quantization, with accuracy degradation observed due to undetected faults when the AUC is not 1. From ATMR(8,4,4), both restorations show enough inference accuracy.

TABLE XI
AUC OF LAYER3.1.1 IN RESNET-20.

Fault_prob[%]	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
ATMR																				
(8,3,3)	0.664	0.796	0.885	0.937	0.967	0.984	0.993	0.997	0.999	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
(8,4,4)	0.838	0.960	0.992	0.999	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
(8,5,5)	0.957	0.999	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
(8,6,6)	0.996	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
(8,7,7)	0.999	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

TABLE XII
RESTORATION BY LINEAR QUANTIZATION FOR LAYER3.1.1 IN RESNET-20.

Fault_prob[%]	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
ATMR																				
(8,3,3)	91.51	91.19	90.70	89.77	89.42	88.61	88.70	88.76	89.71	90.53	90.96	91.16	91.19	91.20	91.20	91.20	91.20	91.20	91.20	91.20
(8,4,4)	91.36	91.05	91.10	91.68	91.82	91.90	91.90	91.90	91.90	91.90	91.90	91.90	91.90	91.90	91.90	91.90	91.90	91.90	91.90	91.90
(8,5,5)	91.77	91.82	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95	91.95
(8,6,6)	91.80	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86	91.86
(8,7,7)	91.91	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89	91.89

TABLE XIII
RESTORATION BY RESIDUAL QUANTIZATION FOR LAYER3.1.1 IN RESNET-20.

Fault_prob[%]	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
ATMR																				
(8,3,3)	91.79	91.03	90.04	89.54	89.00	87.95	87.58	87.78	88.08	89.10	89.28	89.55	89.56	89.58	89.58	89.58	89.58	89.58	89.58	89.58
(8,4,4)	91.38	91.00	90.84	91.31	91.57	91.58	91.60	91.60	91.60	91.60	91.60	91.60	91.60	91.60	91.60	91.60	91.60	91.60	91.60	91.60
(8,5,5)	91.58	91.78	91.84	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82	91.82
(8,6,6)	91.74	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78	91.78
(8,7,7)	91.79	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84	91.84

V. DISCUSSION

From the evaluation of ResNet-20, ATMR of *stem.conv1* cannot effectively handle faults without high bit quantization due to missing errors. In contrast, other layers can achieve near-original accuracy recovery even with 3- or 4-bit quantization. The *stem.conv1* is the convolution layer with the fewest number of filters, and other layers that are placed later gradually increase the number of filters. Therefore, large circuit resources and power savings are expected when all TMR units of ResNet-20 except *stem.conv1* are replaced with ATMR units with 3 or 4 bits. Moreover, using residual quantization for fault detection in *stem.conv1* has the potential to improve the AUC with 3- or 4-bit quantization.

VI. CONCLUSION

In this paper, we propose an ATMR system capable of detecting faults and continuing inference with output restoration, resulting in reduced circuit area and power consumption compared to conventional TMR. Additionally, we propose using output restoration methods based on linear quantization and residual quantization. For evaluation, we conducted experiments that were a simulation of output restoration using ResNet-20 in software. Results from ResNet-20 showed that 3- or 4-bit restoration was effective in layers except for *stem.conv1*. In our future work, we plan to implement ATMR circuits on an FPGA to measure resource utilization and power consumption.

REFERENCES

[1] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, 1998.

[2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

[3] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Quantized neural networks: Training neural networks with low precision weights and activations," *The Journal of Machine Learning Research, Volume 18, Issue 1*, pp 6869–6898, 2016.

[4] R. Banner, I. Hubara, E. Hoffer, and D. Soudry, "Scalable methods for 8-bit training of neural networks," *Neural Information Processing Systems (NIPS)*, 2018.

[5] P. Micikevicius, S. Narang, J. Alben, G. Diamos, E. Elsen, D. Garcia, B. Ginsburg, M. Houston, O. Kuchaiev, G. Venkatesh, and H. Wu, "Mixed precision training," *International Conference on Learning Representations (ICLR)*, 2017.

[6] E. Marica and G. Gielen, "Transistor aging-induced degradation of analog circuits: Impact analysis and design guidelines," *European Solid State Circuits Conference (ESSCIRC)*, 2011.

[7] Y. Saikawa, Y. Owada, Y. Tomioka, H. Saito, and Y. Kohira, "Fault detectable convolutional neural network circuits with dual modular redundancy based on mixed-precision quantization," in *IEICE Technical Report (VLD2023-122)*, 02 2024, pp. 99–140, (In Japanese).

[8] H. Xu, Z. Chen, W. Wu, Z. Jin, S.-y. Kuo, and M. Lyu, "Nv-dnn: Towards fault-tolerant dnn systems with n-version programming," *IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*, 2019.

[9] S. Salamin, G. Zervakis, O. Spantidi, I. Anagnostopoulos, J. Henkel, and H. Amrouch, "Reliability-aware quantization for anti-aging npus," in *2021 Design, Automation and Test in Europe Conference and Exhibition (DATE)*, 2021, pp. 1460–1465.

[10] Y. Li, Y. Liu, M. Li, Y. Tian, B. Luo, and Q. Xu, "D2nn: a fine-grained dual modular redundancy framework for deep neural networks," in *Proceedings of the 35th Annual Computer Security Applications Conference*, 2019, pp. 138–147.

[11] Y. Owada, Y. Tomioka, and H. Saito, "Dual modular redundancy unit of convolutional layer for low-cost and reliable cnns," in *2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 2022, pp. 379–384.

[12] E. Yvinec, A. Dapgony, M. Cord, and K. Bailly, "Rex: Data-free residual quantization error expansion," 2023. [Online]. Available: <https://arxiv.org/abs/2203.14645>