

SoWaF: Shuffling of Weights and Feature Maps: A Novel Hardware Intrinsic Attack (HIA) on Convolutional Neural Network (CNN)

Tolulope A. Odetola*, Syed Rafay Hasan*,

*Department of Electrical and Computer Engineering, Tennessee Technological University, Cookeville, TN 38505, USA

Abstract—Security of inference phase deployment of Convolutional neural network (CNN) into resource constrained embedded systems (e.g. low end FPGAs) is a growing research area. Using secure practices, third party FPGA designers can be provided with no knowledge of initial and final classification layers. In this work, we demonstrate that hardware intrinsic attack (HIA) in such a “secure” design is still possible. Proposed HIA is inserted inside mathematical operations of individual layers of CNN, which propagates erroneous operations in all the subsequent CNN layers that leads to misclassification. The attack is non-periodic and completely random, hence it becomes difficult to detect. Five different attack scenarios with respect to each CNN layer are designed and evaluated based on the overhead resources and the rate of triggering in comparison to the original implementation. Our results for two CNN architectures show that in all the attack scenarios, additional latency is negligible ($< 0.61\%$), increment in DSP, LUT, FF is also less than 2.36% . Three attack scenarios does not require any additional BRAM resources, while in two scenarios BRAM increases, which compensates with the corresponding decrease in FF and LUTs. To the authors’ best knowledge this work is the first to address the hardware intrinsic CNN attack with attacker does not have knowledge of the full CNN.

Index Terms—Convolutional Neural Network, FPGA, Trojan

I. INTRODUCTION

¹ FPGA based Convolutional Neural Network (CNN) inference has gained attention in recent times [1]. FPGA hardware accelerators offer good performance, high energy efficiency, fast prototyping, and capability of reconfiguration, [2], [3]. The re-configurable nature of FPGAs permits flexibility in the mapping of CNNs on FPGA with high accuracy and low precision [4]. The adoption of High Level Synthesis (HLS) in the mapping of CNN on FPGA allows software specifications of accelerators to be synthesizable to hardware [4]. To achieve short time-to-market, the mapping of pre-trained CNN on hardware accelerators is often outsourced to untrusted third parties. They contribute to FPGA design flow, by providing soft IPs or hard IPs (such as bitstream file). Due to their untrusted nature hardware intrinsic security can be compromised via malicious hardware insertions, which are very difficult to detect, especially if the IP is provided as a bitstream file.

Different techniques of inserting hardware attacks into CNNs have been explored. Clements et. al [5] presents a hardware Trojan framework introduced in IP designs. This hardware Trojan generates small bounded perturbations that are added to feature maps of targeted layers of the CNN and causes deterioration in the performance of CNN. Liu et. al in [6] discusses an attack on neural networks where

samples of the input data are generated from the pre-trained model to design a trigger that activates a payload to cause misclassification. These attacks require a manipulation of the CNN parameters to achieve misclassification, which can be detected by carrying out a model integrity test on the hardware design. Moreover, traditionally it is assumed that the attacker has full knowledge of the CNN architecture. We argue that in an effort to deter hardware attacks, the project owner may hide the details of dataset by not providing details of first layer and eliminate the last layers to conceal the classification information as is the case for edge offloading for CNNs [7], [8]. This results in third party IP designers (potential attacker) having no means to evaluate the effectiveness and stealthiness of the attack. Hence, in this paper we demonstrate a framework of attacks called SoWaF (Shuffling of Weights and Feature Maps to corrupt the mathematical computation) that leads to misclassification, without any knowledge of dataset and final classification layer.

A. Motivational Analysis

The above discussed attack scenario lead to the question that, what is the feasibility of hardware intrinsic attack (HIA) if intermittent changes are made to mathematical operations of one of the layers in CNN? The premise is that if a subtle (and stealthy) minimal change in some mathematical operations can lead to misclassification, then it is extremely difficult to detect such attacks. To understand the effect of minimal change in mathematical operations we took a $3 \times 12 \times 12$ input feature map and perform convolution with a channel $3 \times 3 \times 5 \times 5$ weight matrix to obtain a $3 \times 8 \times 8$ output feature maps O_1 . This serve as a baseline result. To device a possible attack, the channels of the weight matrix are then randomly shuffled and used to perform convolution with the input feature map to obtain another $3 \times 8 \times 8$ output feature maps O_2 . Element wise comparison of O_1 and O_2 , shows that 72% of the values are changed more than 95% . This toy example inspired us to do further investigation and see the effect of SoWaF in complete CNN architecture.

B. Research Challenges

We formulated the following research challenges based on the motivational analysis:

- How can an attack be triggered randomly, with the payload activated only intermittently, so that it cannot be detected easily?
- How the malicious changes in the mathematical operations are implemented such that it requires minimal resources but still capable of inducing an effective attack?

¹A version of this work will be published in ISCAS 2021

C. Novel Contribution and Concept Overview

To address the aforementioned research challenges, we propose a HIA methodology for FPGA based CNN inference called SoWaF. The HIA comprises of two stages namely: Offline pre-processing and Runtime payload analysis. The attack is designed to be intermittently triggered. Overview of SoWaF methodology flow is shown in Fig. 1. The red shaded portion of Fig. 1 shows our contribution. Section 1 of the methodology flow involves the offline analysis of the output feature maps to design a stealthy trigger. Section 2 shows the comparison of the additional hardware overhead incurred by the HIA circuitry with the design constraints. Section 3 shows the evaluation of the stealthiness and effectiveness of the attack. Our methodology employs the following analysis and methods:

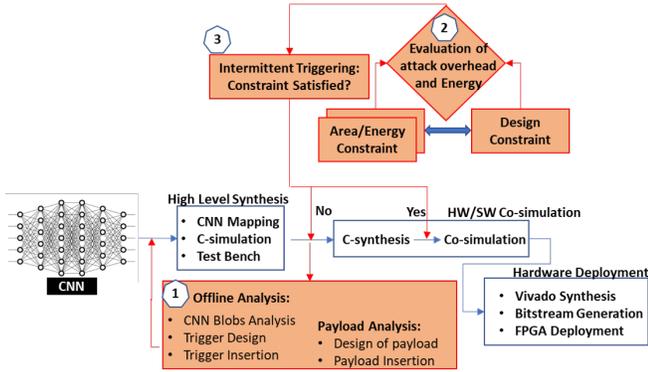


Fig. 1: Design time flow of the HIA methodology. Highlighted boxes represent the novel contributions

- To design a trigger, we propose the exploitation of computation of the layer-by-layer feature maps.
- To achieve stealthiness, we propose a positional relationship based probabilistic trigger, which is intermittent.
- To achieve misclassification with minimum resource overhead, we propose a novel payload that disrupts certain mathematical operations with very minimal (if any) added resources.

The remainder of this paper is organized as follows: Section II describes the threat model. Section III discusses the proposed attack design. Section IV shows experimental results and discussion. Section V provides comparison with state-of-the-art and Section VI concludes the paper.

II. THREAT MODEL

This work proposes a gray-box attack where the attacker has little knowledge of the CNN architecture. We assume that the third party IP designer is not trustworthy. It is assumed that the attacker has no access of the training and testing data samples, i.e. attacker is only designing a CNN without its head (last layers) and initial layers. The attacker provides the implemented CNN hardware design as a bitstream file to the defender (project owner). Fig. 2 shows the trusted and untrusted sections of the design. The 3rd party designer has access to the untrusted sections based on specifications and requirements provided by the trusted party. It is also assumed

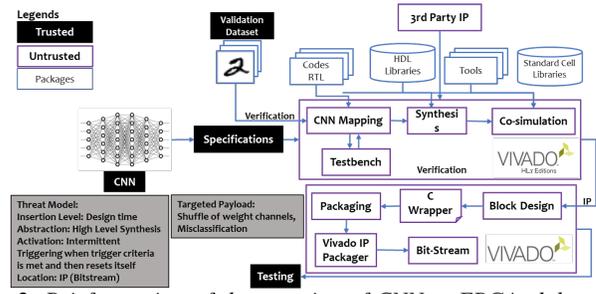


Fig. 2: Brief overview of the mapping of CNN to FPGA elaborating on the threat model and the corresponding payloads.

that for verification purposes, a hardware validation dataset is provided to the 3rd party designer (a normal industrial practice [2]), without revealing any information about the initial layers.

III. SOWAF ATTACK METHODOLOGY

The proposed methodology is sub-divided into 2 phases namely:

- Offline Pre-processing: Trigger Design
- Runtime: Payload Operation

A. Offline Pre-processing: Trigger Design

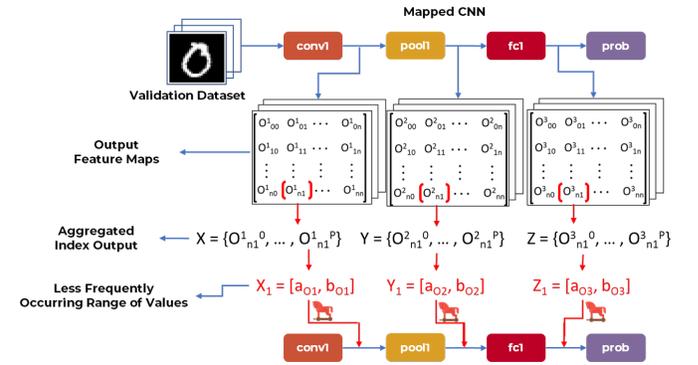


Fig. 3: Offline pre-processing (Trigger Design): Conceptual Representation of offline processing: the aggregation of values from a chosen index helps the selection of less occurring range of values (RoV) (triggers) and the insertion of the trigger (where P is the size of the Validation dataset, O^1 , O^2 , O^3 are the output feature maps of conv1, pool1 and fcl layers, respectively.)

During the functional verification stage the hardware validation dataset can be used by the attacker to access the respective CNN layer's output feature maps for all the dataset. As illustrated in Fig. 3, the attacker assess the statistical properties of the output feature maps to setup a trigger.

In this work, to design the trigger we randomly choose an index ($O^w_{n,m}$, where w is the layer, tuple (n, m) represents the rows and columns of the index) of one of the randomly chosen channels of the output feature map of a targeted layer as shown in line 1 of Algorithm 1. During verification, the attacker may monitor the the values (X or Y or Z) as shown in Fig. 3 of the randomly selected index ($O^w_{n,m}$) against the validation dataset to obtain the range of values (RoV) (where $[a_w, b_w]$ represent the minimum and maximum value of RoV respectively) that are likely to occur at a particular index. This

serves as a sample space to estimate RoV that occur less often on the chosen index as shown in line 3 - 5 of Algorithm 1. From the Aggregated index outputs $\{O_{n,m}^0, \dots, O_{n,m}^P\}$, where $0, \dots, P$ represent each data instance in the validation dataset), as shown in Fig. 3, of a chosen index, we select a RoV whose number of occurrence in the validation dataset satisfy a chosen threshold ($T(O_{n,m}^w) \rightarrow c([a_w, b_w]) = M$) as shown in Line 6 - 10 of Algorithm 1. The selected RoV $([a_w, b_w])$ for a given CNN layer serve as the trigger for the HIA. This offline pre-processing algorithm, Algorithm 1, enables the proposed attack to assess the RoV and to select a stealthy trigger while processing an image.

Algorithm 1 Offline Pre-processing: Trigger Design

Require: Mapping of the CNN to the desired hardware in HLS (C++).

Require: Verification of mapped CNN hardware design.

- 1: Select CNN layer index $O_{n,m}^w$
 - 2: **for** each image $(X) \in$ validation dataset (of size P) **do**
 - 3: $A = \{O_{n,m}^0, \dots, O_{n,m}^P\} \in O_{n,m}^w$
 where: $\{O_{n,m}^0, \dots, O_{n,m}^P\}$ are the numerical values of the chosen index for each data instance in the validation dataset (1, 2, ... P)
 A is output feature map of any layer (X or Y or Z)
 O^w is the chosen channel of the targeted layer
 n, m are the row/column indexes of the chosen channel
 w represents the targeted layer
 - 4: **end for**
 - 5: Select less frequently occurring RoV $[a_w, b_w]$ from A
 - 6: **if** $O_{n,m}^w : T(O_{n,m}^w)$ **then**
 - 7: Select $[a_w, b_w]$ from A
 where: $T(O_{n,m}^w) \rightarrow c([a_w, b_w]) = M$
 c = number of elements in A within $[a_w, b_w]$
 M = Chosen threshold for c
 T is the function that returns a Boolean if the numerical value of the index satisfies the chosen RoV
 $[a_w, b_w]$ are the lower/upper limit of the chosen RoV
 - 8: **else**
 - 9: Select new range $[a_w, b_w]$ and repeat steps 4 to 6.
 - 10: **end if**
 - 11: Insert $[a_w, b_w]$ as trigger in mapped CNN
-

B. Runtime: Payload Operation

To design the payload we proposed an algorithm, Algorithm 2. The payload monitors the selected CNN layer, the selection of CNN layer depends on the additional resource overhead incurred due to the targeted layer and the rate of triggering on the layer. Upon triggering, for convolution layers, the payload shuffles the channels of the weight matrix with another one as illustrated on the right hand side of the decision block in Fig. 4. Line 4 of Algorithm 2 makes sure that the individual channels in a particular layer are swapped, i.e. $Q_w[j_0]$ is swapped with $Q_w[j_f]$ (where $0 < f < l$ and l is the amount of weight matrix channels). Because CNN layers other than convolution and fully connected layers (such as Pooling, etc.)

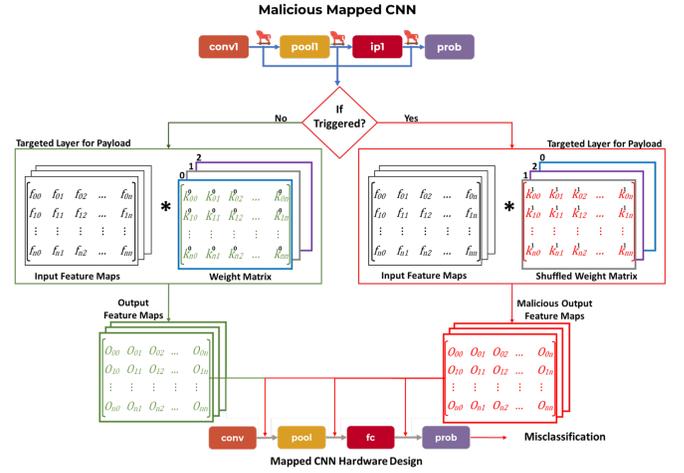


Fig. 4: Representation of the payload during runtime: the right side of the image (shaded in red) shows the weight matrix channels reshuffled to obtain malicious feature maps to achieve misclassification

do not have weight matrices and channels, therefore they are treated separately in line 5 of Algorithm 2. This payload is empirical crafted offline with several experiments (explained in Section IV) to make sure that the malicious mathematical modifications likely results in misclassification.

Algorithm 2 Runtime: Payload Operation

Require: CNN Deployment.

- 1: **for** each image cycle (Im) **do**
 - 2: Monitor selected CNN layer index $O_{n,m}^w$
 - 3: **if** $O_{n,m}^w \in [a_w, b_w]$ **then**
 - 4: $Q_w[j_0, j_1, \dots, j_l] = Q_w[j_f, j_{f+1}, \dots, j_{l-f}]$ (conv layer)
 - 5: $R_w[j_0, j_1, \dots, j_l] = R_w[j_f, j_{f+1}, \dots, j_{l-f}]$ (pool layer)
 where: f is the order factor of shuffling
 $Q_w[j_0, j_1, \dots, j_f]$ is the default weight matrix order
 $R_w[j_0, j_1, \dots, j_f]$ is the default output channel order
 l is the number of weight channels or output channels
 - 6: **if** $0 < j < \frac{l}{2}$ **then**
 - 7: $\frac{l}{2} < f < l$
 - 8: **else**
 - 9: $0 < f < \frac{l}{2}$
 - 10: **end if**
 - 11: **end if**
 - 12: **end for**
-

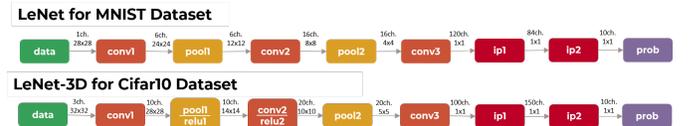


Fig. 5: LeNet and LeNet-3D CNN Models

IV. EXPERIMENT SETUP, RESULTS AND DISCUSSION

The mapped CNN IP is designed using Xilinx's Vivado and Vivado HLS 2018.3 and to generate an IP for resource

TABLE I: Resource overhead comparison between attacks on different layers of LeNet and LeNet-3D compared to their respective originals

Network	Attack Scenario (Sn): Layer	Chs	BRAM	% diff	DSPs	% diff	LUTs (x1000)	% diff	FFs (x1000)	% diff	Latency (x1000) clock-cycles	% diff
LeNet	Original	-	42	-	33	0	118.5	-	58.3	-	680.4	-
	Sn1: conv1 attack	6	42	0	33	0	119.2	+0.61	59.2	+1.5	680.51	+0.003
	Sn2: pool1 attack	6	42	0	33	0	118.9	+0.34	58.8	+0.76	680.51	+0.003
	Sn3: conv2 attack	16	53	+26	33	0	121.3	+2.36	58.8	+0.81	680.58	+0.013
	Sn4: pool2 attack	16	42	0	33	0	119.2	+0.34	59.3	+0.76	680.51	+0.003
	Sn5: conv3 attack	120	162	+285	33	0	780.7	-34	34.5	-41	680.74	+0.038
LeNet-3D for Cifar10	Original	-	59	-	37	-	49.0	-	39.7	-	1685.71	-
	Sn1: conv1 attack	5	59	0	37	0	49.9	+1.81	40.5	+1.8	1685.73	+0.001
	Sn2: pool1 attack	5	59	0	37	0	49.6	+1.16	40.4	+1.76	1685.72	+0.001
	Sn3: conv2 attack	20	79	+34	37	0	48.6	-0.78	39.0	-1.9	1685.72	+0.001
	Sn4: pool2 attack	20	59	0	37	0	50.0	+1.93	41.0	+3.2	1695.99	+0.61
	Sn5: conv3 attack	100	159	+169	37	0	20.1	-59	10.0	-74.6	1685.72	+0.001

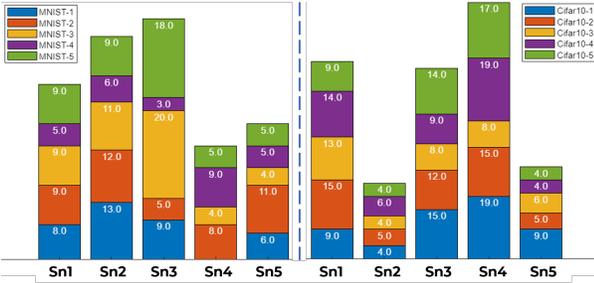


Fig. 6: Graph showing the random nature and low-triggering rate of the different attack scenarios: **Left** represents LeNet trained on MNIST and **Right** and LeNet-3D CNN model

constrained devices. Vivado is used to integrate the generated IP with AXI-interconnects and ZYNQ processor (FPGA ZCU 7020 with clock frequency 100MHz). The HIA is implemented on Lenet (Fig. 5) trained on MNIST dataset and LeNet-3D for Cifar10 datasets, respectively. In this work, we propose 5 different scenarios, where each layer (from conv1 to conv3) is infected with the HIA. Stealthiness (defined as of the additional hardware overhead (BRAM, DSP, flip-flops (FFs), look-up tables (LUTs), Latency) is evaluated for each case and effectiveness (defined as the rate of triggering) of the inserted HIA compared to the original implementation of the mapped CNNs. The attack is carried out on the convolution and pooling layers.

Table I shows that DSP remains same for all the attack scenarios. In both CNNs, BRAM usage remains the same except for Sn3 and Sn5, where BRAM are increased (5th column in Table I). Sn5 compensates this with lower LUTs and FFs usage (see columns 9 and 11 in Table I). For LUTs and FFs in all the scenarios, other than Sn5, (i.e. Sn1 - Sn4) have a very modest increment in usage (up to 2.36%). Similarly, difference in latency between designs with and without HIA design is less than 0.61% in Sn1 - Sn5 (last column in Table I). Hence, we conclude that Sn1, Sn2 and Sn4 can be good choices for an attacker for a stealthy attack, as overall resources and latency effects is minimal. It is observed that for both CNN models, number of weight and output feature map channels are proportional to the additional amount of hardware resources overhead. This can be confirmed from the results of

Sn3 and Sn5 where the higher number of output feature map channels has resulted in higher memory usage. To demonstrate the randomness of the proposed attack, various random input validation dataset is examined. In Fig. 6, for the Sn1, when five sets (200 images each) of data is provided to LeNet and LeNet-3D, the number of trigger occurrences vary randomly between 5 to 9. Same is true for other attack scenarios- making our attack random and stealthy.

V. COMPARISON WITH STATE-OF-THE-ART

We summarized these differences in Table II. Most of the state-of-the-art hardware/firmware attacks on the hardware deployment of CNN requires full knowledge of CNN architecture [5], [7], [8], and [10]. In this paper we argue that to deter the hardware attackers first and last layers may be kept hidden from the un-trusted designers. Hence, our proposed attack is made under more constrained condition. In addition, existing literature requires actual input image for triggering [5], [7] - [10], while proposed design just make use of validation data set. Also in the proposed design, payload implementation does not require extra computation, unlike [5], [7], [8].

TABLE II: Comparison of our approach with other techniques

Criteria	[5]	[9]	[10]	[11]	[12]	Ours
Req. full CNN arch.	✓	✓	✓	x	✓	x
Req. changes in the weights	x	✓	x	x	x	x
Trig. req. Input Image	✓	✓	✓	✓	✓	x
Payload req. extra computation	✓	✓	✓	x	x	x

2

VI. CONCLUSION

To the best of authors' knowledge, this is the first work to propose a HIA targeted at FPGA based CNN inference with attacker having no knowledge of initial layers, datasets, and final classification layer. The attack achieves misclassification by shuffling the weight matrices of convolution layers to propagate wrong feature maps. This attack is carried out without changes in the model parameters. Our results for two CNN architectures show that in all the attack scenarios, additional latency is negligible (< 0.61%), increment in DSP, LUT, FF is also less than 2.36%. Three of the five investigated

²Several works have addressed security and privacy in many applications [13]-[38]

scenarios show very minimal changes in BRAM. Proposed attack is triggered intermittently and our results show that the number of triggers and its occurrence instance are completely random.

REFERENCES

- [1] K. Abdelouahab, M. Pelcat, J. Serot, and F. Berry, "Accelerating cnn inference on fpgas: A survey," *arXiv preprint arXiv:1806.01683*, 2018.
- [2] T. A. Odetola, K. M. Groves, and S. R. Hasan, "2l-3w: 2-level 3-way hardware-software co-verification for the mapping of deep learning architecture (dla) onto fpga boards," *arXiv preprint arXiv:1911.05944*, 2019.
- [3] M. T. Hailesellase and S. R. Hasan, "Mulnet: A flexible cnn processor with higher resource utilization efficiency for constrained devices," *IEEE Access*, vol. 7, pp. 47 509–47 524, 2019.
- [4] J. H. Kim, B. Grady, R. Lian, J. Brothers, and J. H. Anderson, "Fpga-based cnn inference accelerator synthesized from multi-threaded c software," in *2017 30th IEEE International System-on-Chip Conference (SOCC)*. IEEE, 2017, pp. 268–273.
- [5] J. Clements and Y. Lao, "Hardware trojan attacks on neural networks," *arXiv preprint arXiv:1806.05768*, 2018.
- [6] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang, "Trojaning attack on neural networks," 2017.
- [7] R. Hadidi, J. Cao, M. Woodward, M. S. Ryoo, and H. Kim, "Musical chair: Efficient real-time recognition using collaborative iot devices," *arXiv preprint arXiv:1802.02138*, 2018.
- [8] L. Zeng, X. Chen, Z. Zhou, L. Yang, and J. Zhang, "Coedge: Cooperative dnn inference with adaptive workload partitioning over heterogeneous edge devices," *IEEE/ACM Transactions on Networking*, 2020.
- [9] M. Zou, Y. Shi, C. Wang, F. Li, W. Song, and Y. Wang, "Potrojan: powerful neural-level trojan designs in deep learning models," *arXiv preprint arXiv:1802.03043*, 2018.
- [10] J. Clements and Y. Lao, "Hardware trojan design on neural networks," in *2019 IEEE International Symposium on Circuits and Systems (ISCAS)*. IEEE, 2019, pp. 1–5.
- [11] Y. Zhao, X. Hu, S. Li, J. Ye, L. Deng, Y. Ji, J. Xu, D. Wu, and Y. Xie, "Memory trojan attack on neural network accelerators," in *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2019, pp. 1415–1420.
- [12] Z. Liu, J. Ye, X. Hu, H. Li, X. Li, and Y. Hu, "Sequence triggered hardware trojan in neural network accelerator," in *2020 IEEE 38th VLSI Test Symposium (VTS)*. IEEE, 2020, pp. 1–6.