

# Drop-Connect as a Fault-Tolerance Approach for RRAM-based Deep Neural Network Accelerators

Mingyuan Xiang<sup>†\*</sup>, Xuhan Xie<sup>†\*</sup>, Pedro Savarese<sup>‡</sup>, Xin Yuan<sup>†</sup>, Michael Maire<sup>†</sup> and Yanjing Li<sup>†</sup>

<sup>†</sup>Department of Computer Science, University of Chicago, Chicago, IL USA, 60637

<sup>‡</sup>Toyota Technological Institute at Chicago, Chicago, IL USA, 60637

**Abstract**—Resistive random-access memory (RRAM) is widely recognized as a promising emerging hardware platform for deep neural networks (DNNs). Yet, due to manufacturing limitations, current RRAM devices are highly susceptible to hardware defects, which poses a significant challenge to their practical applicability. In this paper, we present a machine learning technique that enables the deployment of defect-prone RRAM accelerators for DNN applications, without necessitating modifying the hardware, retraining of the neural network, or implementing additional detection circuitry/logic. The key idea involves incorporating a drop-connect inspired approach during the training phase of a DNN, where random subsets of weights are selected to emulate fault effects (e.g., set to zero to mimic stuck-at-1 faults), thereby equipping the DNN with the ability to learn and adapt to RRAM defects with the corresponding fault rates. Our results demonstrate the viability of the drop-connect approach, coupled with various algorithm and system-level design and trade-off considerations. We show that, even in the presence of high defect rates (e.g., up to 30%), the degradation of DNN accuracy can be as low as less than 1% compared to that of the fault-free version, while incurring minimal system-level runtime/energy costs.

**Index Terms**—fault tolerance, neural network, machine learning, RRAM

## I. INTRODUCTION

Recent advancements in Deep Neural Networks (DNNs) have demonstrated significant success across various applications. However, the increasing complexity and capabilities of DNNs necessitate substantial computational power and memory bandwidth in conventional Von Neumann architectures to accelerate DNN applications. A promising alternative lies in the utilization of novel architectures constructed with emerging technologies. Among the various options, the Resistive RAM (RRAM) crossbar-based architecture, comprised of memristor cells [1], emerges as an innovative compute-in-memory solution that not only reduces power consumption, but also boosts processing speeds. Illustrated in Fig. 1 is a standard RRAM crossbar. Within the crossbar, DNN kernels are unfolded and embedded with each memristor cells, each of which retaining a single weight value, while input data is continuously streamed into the crossbar from its wordlines. The analog nature of this architecture makes it well-suited for vector-matrix multiplication, as the dot product operation can be replicated using Kirchhoff’s circuit law.

While RRAM-based DNN accelerators offer dramatic improvements in energy efficiency and throughput over tradi-

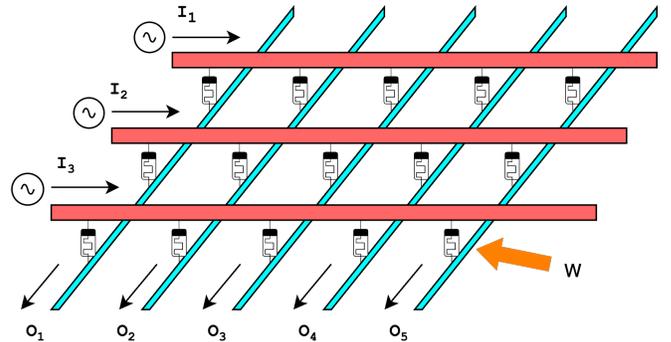


Fig. 1: Example of a RRAM crossbar.

tional architectures, memristor devices are prone to various hardware defects [2], primarily due to manufacturing constraints. The most prominent among these defects manifest as the Stuck-At-Faults (SAF). In particular, the Stuck-At-One (SA1) faults, where memristor cells are trapped in a high state, can lead to both write and read failures so that weight values stored in these faulty cells are always interpreted as 0’s. SA1 faults account for a significant portion (9.04%) of all faults in RRAM crossbars [2].

Even if only a fraction of these faults occur, the impact on DNN model accuracy can be significant. Existing work [3]–[5] typically relies on the actual defect distributions that can be obtained through memory testing [2], [6] to design fault-tolerant RRAM accelerators. However, these techniques either compromise network accuracy (e.g., as seen in remap techniques [4], especially with higher defect rates) or impose substantial deployment overhead (e.g., in retrain-based approaches, where the entire model must be trained from scratch every time it is deployed to a new RRAM crossbar). Moreover, because of the inherent variability and drift associated with the forming and switching operations in memristors, the distribution of defects in a crossbar can evolve over time, which adds another layer of complexity and makes it even more challenging and less attractive to deploy previously proposed solutions. Checksum-based techniques [7] inspired by algorithm-based fault-tolerance are effective and efficient for RRAM DNN fault tolerance as well. However, these techniques require additional hardware resources to compute and compare checksums, which are not always feasible or desirable.

\* These two authors contribute equally.

This paper explores a training method inspired by a machine learning training technique called drop-connect. In drop-connect [8], randomly selected sets of weights of a DNN are set to 0 (and not updated) during different training iterations. While this technique was originally proposed to avoid overfitting during training, it can also be used to emulate random SA1 faults for a given fault rate. This approach is a natural solution to enable RRAM crossbars to tolerate SA1 faults. The intuition is that, if the network can be trained to recognize and understand various fault effects and patterns, taking into account the corresponding fault distributions, then it should be well-equipped to handle similar fault effects and distributions during the inference phase.

The most notable benefits of our approach is that it does not require modifications to the hardware, retraining of the neural network, or the implementation of additional detection circuitry/logic, such as checksum logic. Moreover, this approach offers the opportunity to train the network so that it can adapt to the evolving defect distributions overtime, further enhancing its effectiveness and applicability.

The major contributions of our work are:

- We conduct a thorough investigation and analysis of the drop-connect inspired technique as a fault-tolerance solution for RRAM-based DNN accelerators. More importantly, our technique requires **no information** of actual defect distribution at deployment. The key differences between our work and others are summarized in Table. I
- We provide the detailed implementation of this approach.
- We perform experiments across various representative DNNs and fault rates, and explore various tradeoffs between network accuracy and efficiency (in terms of system-level runtime and energy consumption) to demonstrate the efficacy of our approach.

TABLE I: Comparison of major fault-tolerant techniques for RRAM crossbars.

method	Additional Circuits	Retraining and Mapping
Liu et al. [5]	Defect Detection	required
Das et al. [7]	Checksums	None
Li et al. [9]	Refresh and Detection	required
Ours	<b>None</b>	<b>None</b>

Our key findings and results are summarized below.

- 1) The drop-connect-inspired approach is a viable solution for enabling fault-tolerant RRAM DNN accelerators, particularly when the fault rate is relatively low or when a modest accuracy loss (e.g., 5%) is acceptable. For certain networks (e.g., MobileNet V2), the degradation in accuracy is less than 1% even when the fault rate is very high (30%), while incurring no additional system-level costs.
- 2) The best accuracy levels are achieved when a higher drop-connect rate is used, compared to the expected RRAM fault rate. However, high drop-connect rates (beyond 30% – 40%) adversely affect network accuracy even in fault-free scenarios.

- 3) Further improving accuracy in RRAM crossbars with drop-connect-inspired fault tolerance may involve widening the original network (increasing the number of channels). This compensates for information loss due to weights forced to certain values, but at the cost of higher runtime and energy consumption, and diminishing returns. Systematically exploring trade-offs between system-level costs and network accuracy is therefore crucial. Our results show that, a 20%/60% increase in the number of channels yields up to 4%/12.5% improvements in test accuracy, respectively, compared to 0% increase, while incurring up to 42.6%/153.3% performance/energy costs.
- 4) Due to the unique properties of convolution layers with 1x1 kernels, running these layers in traditional architectures such as CPUs achieve allows the network to achieve higher network and runtime/energy efficiency simultaneously.
- 5) In certain networks, modifying the structure of a few critical layers serves as an alternative approach to accompany drop-connect in order to improve DNN accuracy. For example, in ResNet20 [10], increasing the kernel size in the shortcut layers from 1x1 to 3x3 and deploying these layers in RRAM crossbars achieves comparable network accuracy to deploying the original network layers in fault-free devices.

Our approach is orthogonal to other fault-tolerance techniques such as the post-training remapping strategies, and they can be combined to further enhance DNN accuracy. Moreover, our work has brought to light an important revelation – there exist nuances in the adaptation of machine learning approaches to tackle system-level challenges, which necessitates an in-depth understanding not only of the original machine learning technique, but also of the broader system-level implications and tradeoffs associated with applying an existing method for a new purpose. For example, based on key result III-B discussed above, the “sweet-spot” for the drop-connect rates must be carefully selected to obtain the optimal accuracy, while employing the vanilla drop-connect technique where the drop-connect rate is set to SA1 fault rate is likely to be sub-optimal. This revelation holds broader relevance beyond our specific study, extending to other system/design challenges where machine learning techniques are adopted and adapted as solutions, for which there is an imperative need for meticulous consideration and thorough analysis.

The rest of the paper is structured as follows: Sec. II describes our methodology. Sec. III presents our experimental results, followed by related work in Sec. IV and conclusions in Sec. V.

## II. METHODOLOGY

### A. The Original Drop-Connect Approach

Drop-connect [8] is a regularization technique to prevent DNN models from overfitting during training. A convolution

operation integrated with drop-connect can be defined as:

$$O_{mij}^{(l)} = \frac{1}{1-p} \sum_n \sum_k \sum_s I_{n,i+k,j+s}^{(l)} \cdot W_{mnks}^{(l)} \cdot M_{mnks}^{(l)} \quad (1)$$

Here,  $I_{nij}^{(l)}$ ,  $O_{mij}^{(l)}$ , and  $W_{mnks}^{(l)}$  are inputs, outputs, and weights of the convolution layer  $l$ .  $M^{(l)}$  is a mask that satisfies the Bernoulli distribution, i.e.,  $M^{(l)} \sim \text{Bernoulli}(p)$ .  $p$  is the drop-connect probability. Additionally, the output needs to upscale by  $\frac{1}{1-p}$  during training to maintain the expected distribution.

For networks that incorporate normalization layers, such as batch normalization [11], there is an important implementation detail that is worth noting. A scaling factor of  $\frac{1}{1-p}$  must also be applied to normalize the weights when drop-connect is applied [8]. This ensures that only non-zero weights are normalized to preserve the correct weight value distribution.

### B. Drop-Connect Adapted for Fault-Tolerance Purposes

As discussed in Sec. I, drop-connect provides a way for DNN models to learn to compensate for RRAM defects regardless of the actual defect cell distribution while maintaining high accuracy. As such, no additional costs and overheads are required during deployment time.

An important question for adapting drop-connect for fault-tolerance purposes is the following: what is the optimal drop-connect rate during training? To answer this question, we sweep the drop connect rate for a given fault rate (see Sec. III). In these experiments, as the running statistics (i.e., moving mean and variance) of batch normalization can differ between the drop-connect rate and the actual defect rate, our training algorithm involves an additional epoch to align the scaling factor of normalization layers with the SA1 fault rate during inference (which is different from the drop-connect rate during training). As shown in Algorithm 1, during this epoch, we apply the inference-time scaling factor (i.e., RRAM SA1 fault rate) to adjust the moving mean/variance values, matching them more closely to the actual SA1 fault distribution, while keeping the network weights constant.

---

#### Algorithm 1 UpdateVar(MODEL)

---

```

1: //  $p'$ : inference-time scaling factor, i.e., SA1 fault rate
2: for  $p'$  in [0%, 10%, 20%, 30%] do
3:   MODEL.train()
4:   for  $batch$  in TrainingSet do
5:     Freeze MODEL.weights
6:     MODEL.dropConnect( $p'$ )
7:     MODEL.forward( $batch$ )
8:   end for
9: end for

```

---

Moreover, we explore widening the original network (i.e., increasing the number of channels in convolution layers) to compensate for information loss due to drop-connect, as shown in Fig. 2 (highlighted in red).  $p$  is the percentage increase in the number of channels.

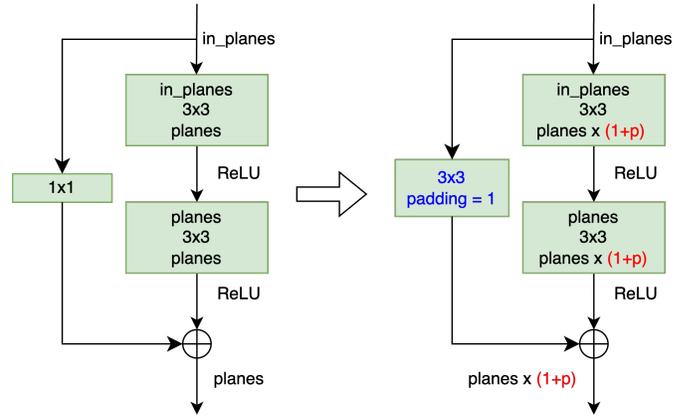


Fig. 2: Increasing network width (red) and/or increasing kernel size of 1x1 convolution layers (blue) to compensate for information loss due to drop-connect.

### C. Special Considerations on Convolution Layers with 1x1 Kernels

When applying the drop-connect approach during training, it is not applied to convolution layers where the kernel size is 1x1. Such 1x1 convolution layers are prevalent in modern DNN models, such as point-wise convolution in MobileNet V2 [12], or shortcut convolution in ResNet20 [10]. Unlike convolutions with kernel sizes larger than 1x1, these 1x1 convolution layers can be efficiently optimized on traditional processors, such as CPUs. Specifically, 1x1 convolutions are equivalent to multiplication and accumulation (MAC) operations along the input channel dimensions. These MAC operations can be readily computed using optimized architecture extensions such as instructions in the Intel Advanced Vector Extensions (AVX) [13]. Furthermore, loop unrolling can be employed to mitigate the loop branching overhead. Consequently, 1x1 convolution layers can be executed efficiently on traditional processors, so it is not critical to offload them to accelerators such as RRAM crossbars.

Coincidentally, these 1x1 convolution layers also constitute a critical component of the network model. For instance, the point-wise convolution in MobileNet V2 increases the output channels for the separable convolution. When an SA1 fault occurs, the corresponding output channels are all zeroed, leading to a significant loss of information. At the same time, drop-connect does not work well in these layers, again due to the significant loss of information (results in Sec. III).

In summary, for networks with 1x1 convolution layers, considering both system-level runtime/energy efficiency and network accuracy, it is more advantageous to execute these 1x1 convolution layers on traditional fault-free devices, such as CPUs and GPUs.

### D. Modifying Network Layer Structure to Enhance Drop-Connect

Motivated by the promising application of the drop-connect approach for fault-tolerance and the crucial role observed in

1x1 convolution layers for neural network robustness, we explore an alternative machine learning technique which involves modifying the structure of the 1x1 convolution layers. The idea is to increase the kernel size – for example, from 1x1 to 3x3. In the original network, the drop-connect-inspired approach encounters challenges with the 1x1 convolution layers, primarily stemming from notable information loss within these layers. By expanding the kernel size, we can now apply drop-connect to these layers, enabling their deployment on RRAM-based DNN accelerators. Fig 2 provides an example of how the 1x1 shortcut layer is modified (highlight in blue). In this example, each side of the input is padded with 0, and the 1x1 convolution is substituted with a 3x3 convolution. Note that, these 3x3 shortcut layers serve the same role as identity mapping in residual neural networks.

### III. EVALUATION AND KEY RESULTS

#### A. Experiment Setup

We demonstrate the effectiveness of our approach through simulation experiments. We simulate the RRAM crossbar behaviors and fault effects using an in-house PyTorch-based [14] simulator. In our simulation, each memristor cell stores an 8-bit weight value, a common RRAM crossbar configuration [15]–[17].

We included a set of representative DNN benchmarks in our experiments, including ResNet20 [10], MobileNet V2 [12] and VGG13 [18] on the CIFAR-10 dataset [19]. Given a fault rate, after training these networks using the drop-connect-based approach discussed in Sec. II, we construct 100 RRAM crossbars with random faults following the given fault rate, and apply the trained model on each of the crossbars. We collect and report the test accuracy by taking the average of these 100 runs.

We performed a systematic analysis by sweeping fault rates and drop-connect rates, increasing the number of channels, and experimenting with expanding the kernel size of short-cut layers in ResNet20 from 1x1 to 3x3.

#### B. Results for Adapting Drop-Connect as a Fault-Tolerance Solution

The accuracy results for different drop-connect and fault rate combinations are shown in Figs. 3, 4, and 5 for different networks.

Two key observations can be made from these results. First, the drop-connect-based approach is effective in enabling fault-tolerant RRAM-based DNN accelerators, when the SA1 fault rate is up to 10%. For example, for a 10% fault rate, the degradation in network accuracy is less than 2% for VGG and MobileNet V2. For higher fault rates, larger gaps in accuracy are seen – 4-10%/6-20% for a 20%/30% fault rate, respectively. However, this approach still provides substantial benefit in model accuracy compared to networks training without drop-connect (i.e., the data points correspond to 0% drop-connect rate).

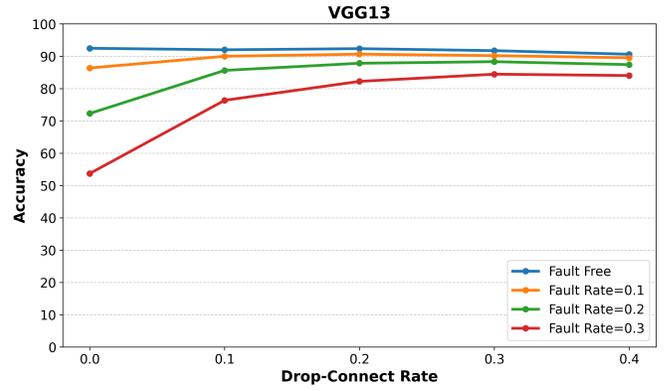


Fig. 3: Network Accuracy of VGG13 for different drop-connect and fault rates.

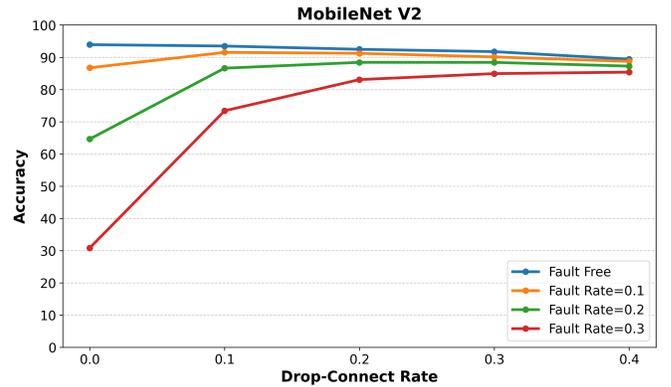


Fig. 4: Network Accuracy of MobileNet V2 for different drop-connect and fault rates.

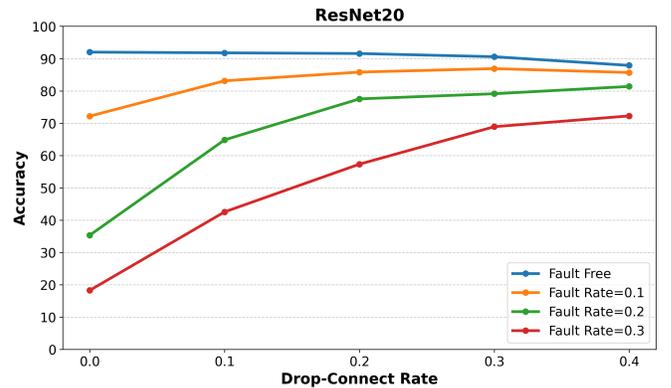


Fig. 5: Network Accuracy of ResNet20 for different drop-connect and fault rates

Second, the best accuracy for each fault-rate is achieved when the drop-connect rate is higher. The best drop-connect rate for a given fault-rate is shown in Figs. 6a and 6b.

Overall, our approach works well for VGG, since it is an over-parameterized network with inherently higher fault-tolerance than other networks. For MobileNet V2, our approach achieves similar levels of effectiveness. This may be attributed to executing the critical point-wise layers in fault-

free hardware. On the other hand, for ResNet20 which is a more challenging network, the drop-connect approach alone may not be able to achieve the desired level of accuracy.

Network accuracy improves as we increase the width of the network, as shown in Figs. 6a and 6b. Note that, we did not run the same experiments for MobileNet V2 because the accuracy is already very high even with the original network configuration.

More channels come with more filters in convolution layers, allowing the model to learn a greater number of features from the data. In addition, having more channels provides the model with more flexibility to determine which filter is crucial for the task, which leads to higher accuracy. This method is particularly pronounced for ResNet20. Its impact becomes more evident when evaluating the model on RRAM crossbars with higher fault rates. With a 10%/20% fault rate, the degradation in accuracy compared to the fault-free ResNet20 is only 0.2%/3%, respectively. VGG also benefits from this method, achieving at least 90% accuracy (< 2% degradation) on RRAM crossbars with fault rates ranging from 0 to 30%.

However, there is a trade-off here. Increasing the number of channels inevitably leads to higher computational costs and longer runtime. We estimate the energy consumption and latency with respect to the width of networks, using the computational efficiency and power efficiency results from the 64-chip ISAAC RRAM accelerator [1], and the results are shown in Table. II. The width of each network is normalized to that of the original model (first row in the table).

TABLE II: Latency and energy estimation for different network configurations with different network widths (normalized to the original network width).

Width of Network	ResNet20	MobileNet V2	VGG13
	latency (s), energy (kW)		
1x	15.74	0.68	87.83
	113.32	4.89	632.16
1.2x	22.13	0.96	125.25
	159.28	6.93	901.56
1.4x	29.80	—	170.59
	214.53	—	1227.84
1.6x	39.11	—	222.53
	281.56	—	1601.71

Increasing network width can also lead to potential overfitting problems as well as diminishing returns. For instance, we observe that the benefit of increasing network width for VGG starts to diminish at 40% increase in the number of channels, and increasing the network width further does not significantly improve accuracy. Moreover, for ResNet20, with a fault rate of 30%, the accuracy degradation of around 7% may still be too high, even though it is a substantial improvement over the network accuracy achieved using the original network configuration (the accuracy degradation is around 79%).

Therefore, given the various tradeoffs associated with this approach, it is crucial to systematically explore and select the best design point that balances network accuracy and

runtime/energy costs for a given use scenario. It is also worth noting that, even with higher system-level costs, our approach can be still more desirable than others that require additional hardware support (e.g., special RRAM circuit or peripheral checksum logic), which requires more design efforts and cannot be readily deployed on existing hardware.

### C. Results for Increasing the Size of Kernels

In Fig. 6c, we show the accuracy for ResNet20 when the kernel size of the short-cut layers is increased from 1x1 to 3x3, with drop-connect applied to these layers during training. The results are comparable to those shown in Fig. 6b, suggesting that this is a promising alternative approach to achieve fault tolerance.

### D. Results Demonstrating the Fault Criticality of Convolution Layers with 1x1 Kernels

In Fig. 7, we demonstrate the fault criticality of convolution layers with 1x1 kernels by comparing the difference between: (1) our approach, which is to execute these layers in traditional fault-free architectures, and (2) applying drop-connect and mapping these layers to faulty RRAM crossbars. The dramatic difference is clear evidence that executing these layers in a fault-free manner is crucial to achieve high network accuracy.

## IV. RELATED WORK

RRAM defects and faults have been extensively examined and characterized in previous work [2], [20], [21]. They can be divided into two main categories: soft errors [20] and hard errors [2]. In the scope of this paper, the focus is on hard errors, because solutions at the circuit level have already been devised to address soft RRAM errors [20]. Hard errors include the stuck-at and transition faults. Although traditional methods such as the March-C algorithm [6] or the squeeze-search algorithm [2] are effective in detecting hard RRAM errors, they cannot mitigate the defects [22]. Thus, fault-tolerance techniques are required for practical deployment of RRAM accelerators.

A retrain and post-mapping method is proposed by [5], where they detect the defect distribution and retrain the DNN models. After training, a remapping scheme is adopted such that the least important weights are mapped to faulty memristor cells. Similar retrain techniques are also proposed in [3], [4], [9], [23]. A significant limitation of retrain-based techniques is that they require training an entire neural network from scratch every time the network is deployed to a new accelerator (or when the fault distribution of an accelerator changes), which incurs prohibitively high resource overheads, and may not even be possible because the training dataset may not be available. Post-mapping methods avoid the high overhead of retrain, but they alone cannot achieve acceptable levels of network accuracy.

Multiple circuit-level solutions have been proposed to handle RRAM defects. For example, [7] employs a parity matrix for a majority-vote-based checksum across the entire crossbar.

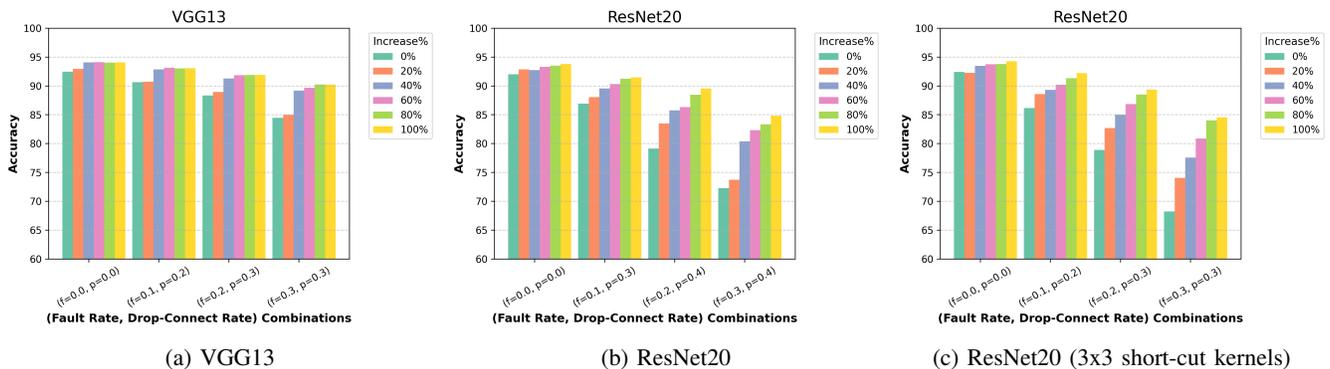


Fig. 6: (a) and (b) the combinations of the fault rate and the drop-connect rate that achieve the highest network accuracy for different network width; (c) increasing the size of kernels in short-cut layers from 1x1 to 3x3.

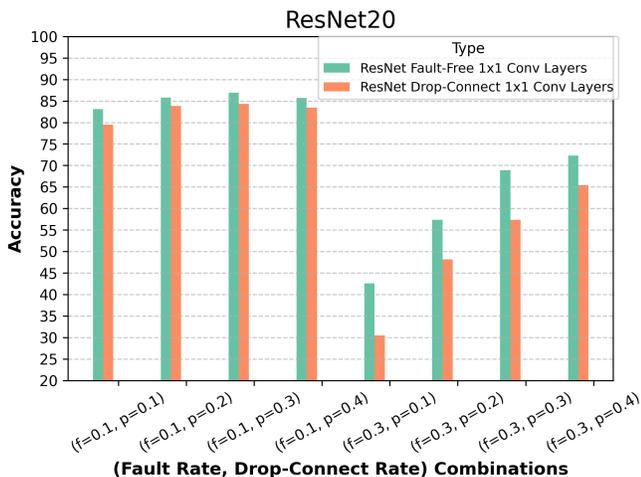


Fig. 7: ResNet20, comparison of fault-free 1x1 convolution layers and when drop-connect and SA1 faults are applied to the same layers.

And in [24], an approach based on the diagonal error correction code is proposed. However, such circuit solutions are only effective when the defect rate is no higher than 5% [25]. For higher defect rates, they introduce non-negligible latency overheads (26% [24]) or memory overheads ranging from 9% to 30% [7].

The concept of drop-connect was mentioned in [26] to achieve fault-tolerance in RRAM crossbars; however, the authors approached its application naively, lacking the profound understanding and comprehensive system-level trade-off analysis that our work provides — a necessity for the successful adaptation of this technique. Consequently, their work lacks meaningful insights and falls short in delivering desirable results. Notably, results are only reported for a limited set of simple and shallow models (e.g., LeNet [27] and AlexNet [28]), yet acceptable levels of network accuracy still cannot be achieved – the training accuracy of MNIST is 78.34%, a significant degradation vs. the fault-free case where a 92.8%

training accuracy can be achieved. A pruning-based scheme was proposed by [29] under the intuition that defects can be mitigated if the positions of pruned weights overlapped with faulty memristor cells. However, this method also suffers from significant network accuracy degradation with high fault rates. For example, the accuracy of ResNet18 using the CIFAR-10 dataset drops to less than 20% even for a low fault rate of 10% in [29].

Drop-connect has been proposed to mitigate the effects of hardware defects on DNNs in systolic array-based DNN accelerators [30]. Moreover, a dropout-based method, where neurons are dropped (instead of weights being dropped in the case of drop-connect) is proposed to overcome faults in the neurons of spiking neural networks [31]. In contrast, our work focuses on RRAM-based DNN accelerators.

## V. CONCLUSION

In this paper, we perform a thorough study on a drop-connect-inspired technique to enable fault tolerance in RRAM-based DNN accelerators. Distinct from previous work, we incorporate various algorithm-/system-level considerations and analysis, and also conduct comprehensive experiments. Through our study, we have obtained various new insights, and the main conclusion is that our approach is viable as a fault-tolerance solution, especially if the fault rate is low and/or if modest network accuracy degradation is acceptable. Our approach allows various tradeoffs between network accuracy and system-level runtime/energy efficiency to be obtained, so that the best design point can be chosen for the specific use scenarios. At the same time, this approach does not require modifications to the hardware, retraining of the neural network, or the implementation of additional detection circuitry/logic. However, in order to tolerate a higher fault rate or close the accuracy gap, other machine learning and system techniques are needed. We plan to build on top of this work and investigate even more efficient fault-tolerance techniques targeting RRAM-based DNN accelerators.

## REFERENCES

- [1] A. Shafiee, A. Nag, N. Muralimanohar, R. Balasubramonian, J. P. Strachan, M. Hu, R. S. Williams, and V. Srikumar, "ISAAC: A con-

- volutional neural network accelerator with in-situ analog arithmetic in crossbars,” in *43rd ACM/IEEE Annual International Symposium on Computer Architecture, ISCA 2016, Seoul, South Korea, June 18-22, 2016*. IEEE Computer Society, 2016, pp. 14–26.
- [2] C.-Y. Chen, H.-C. Shih, C.-W. Wu, C.-H. Lin, P.-F. Chiu, S.-S. Sheu, and F. T. Chen, “Rram defect modeling and failure analysis based on march test and a novel squeeze-search scheme,” *IEEE Transactions on Computers*, vol. 64, no. 1, pp. 180–190, 2015.
  - [3] Y. Huang, Y. He, J. Wang, J. Yue, L. Zhang, K. Zou, H. Yang, and Y. Liu, “Bit-aware fault-tolerant hybrid retraining and remapping schemes for rram-based computing-in-memory systems,” *IEEE Trans. Circuits Syst. II Express Briefs*, vol. 69, no. 7, pp. 3144–3148, 2022.
  - [4] L. Chen, J. Li, Y. Chen, Q. Deng, J. Shen, X. Liang, and L. Jiang, “Accelerator-friendly neural-network training: Learning variations and defects in RRAM crossbar,” in *Design, Automation & Test in Europe Conference & Exhibition, DATE 2017, Lausanne, Switzerland, March 27-31, 2017*, D. Atienza and G. D. Natale, Eds. IEEE, 2017, pp. 19–24.
  - [5] C. Liu, M. Hu, J. P. Strachan, and H. H. Li, “Rescuing memristor-based neuromorphic design with high defects,” in *Proceedings of the 54th Annual Design Automation Conference, DAC 2017, Austin, TX, USA, June 18-22, 2017*. ACM, 2017, pp. 87:1–87:6.
  - [6] Y.-X. Chen and J.-F. Li, “Fault modeling and testing of 1T1r memristor memories,” in *2015 IEEE 33rd VLSI Test Symposium (VTS)*, 2015.
  - [7] A. Das and N. A. Toubba, “Selective checksum based on-line error correction for RRAM based matrix operations,” in *38th IEEE VLSI Test Symposium, VTS 2020, San Diego, CA, USA, April 5-8, 2020*. IEEE, 2020, pp. 1–6.
  - [8] L. Wan, M. D. Zeiler, S. Zhang, Y. LeCun, and R. Fergus, “Regularization of neural networks using dropconnect,” in *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, ser. JMLR Workshop and Conference Proceedings, vol. 28. JMLR.org, 2013, pp. 1058–1066.
  - [9] W. Li, Y. Wang, H. Li, and X. Li, “Rramedy: Protecting rram-based neural network from permanent and soft faults during its lifetime,” in *37th IEEE International Conference on Computer Design, ICCD 2019, Abu Dhabi, United Arab Emirates, November 17-20, 2019*. IEEE, 2019, pp. 91–99. [Online]. Available: <https://doi.org/10.1109/ICCD46524.2019.00020>
  - [10] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. IEEE Computer Society, 2016, pp. 770–778.
  - [11] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, ser. JMLR Workshop and Conference Proceedings, F. R. Bach and D. M. Blei, Eds., vol. 37. JMLR.org, 2015, pp. 448–456.
  - [12] M. Sandler, A. G. Howard, M. Zhu, A. Zhmoginov, and L. Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. Computer Vision Foundation / IEEE Computer Society, 2018, pp. 4510–4520.
  - [13] *Intel Intrinsics Guide*, 2023. [Online]. Available: <https://www.intel.com/content/www/us/en/docs/intrinsics-guide/index.html>
  - [14] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Z. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 8024–8035.
  - [15] P. Yao, H. Wu, B. Gao, J. Tang, Q. Zhang, W. Zhang, J. J. Yang, and H. Qian, “Fully hardware-implemented memristor convolutional neural network,” *Nature*, vol. 577, no. 7792, pp. 641–646, 2020.
  - [16] Y. Wang, T. Tang, L. Xia, B. Li, P. Gu, H. Yang, H. Li, and Y. Xie, “Energy efficient RRAM spiking neural network for real time classification,” in *Proceedings of the 25th edition on Great Lakes Symposium on VLSI, GLVLSI 2015, Pittsburgh, PA, USA, May 20 - 22, 2015*, A. K. Jones, H. H. Li, A. K. Coskun, and M. Margala, Eds. ACM, 2015, pp. 189–194.
  - [17] Z. He, J. Lin, R. Ewetz, J. Yuan, and D. Fan, “Noise injection adaption: End-to-end rram crossbar non-ideal effect adaption for neural network mapping,” in *Proceedings of the 56th Annual Design Automation Conference 2019, DAC 2019, Las Vegas, NV, USA, June 02-06, 2019*. ACM, 2019, p. 57.
  - [18] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.
  - [19] A. Krizhevsky, V. Nair, and G. Hinton, “Cifar-10 (canadian institute for advanced research).”
  - [20] A. M. S. Tossou, M. H. Anis, and L. Wei, “RRAM refresh circuit: A proposed solution to resolve the soft-error failures for hfo<sub>2</sub>/hf 1t1r RRAM memory cell,” in *Proceedings of the 26th edition on Great Lakes Symposium on VLSI, GLVLSI 2016, Boston, MA, USA, May 18-20, 2016*, A. K. Coskun, M. Margala, L. Behjat, and J. Han, Eds. ACM, 2016, pp. 227–232.
  - [21] S. Kannan, J. Rajendran, R. Karri, and O. Sinanoglu, “Sneak-path testing of crossbar-based nonvolatile random access memories,” *IEEE Transactions on Nanotechnology*, 2013.
  - [22] A. Chaudhuri and K. Chakrabarty, “Analysis of process variations, defects, and design-induced coupling in memristors,” in *IEEE International Test Conference, ITC 2018, Phoenix, AZ, USA, October 29 - Nov. 1, 2018*. IEEE, 2018, pp. 1–10.
  - [23] L. Xia, M. Liu, X. Ning, K. Chakrabarty, and Y. Wang, “Fault-tolerant training with on-line fault detection for rram-based neural computing systems,” in *Proceedings of the 54th Annual Design Automation Conference, DAC 2017, Austin, TX, USA, June 18-22, 2017*. ACM, 2017, pp. 33:1–33:6.
  - [24] O. Leitersdorf, B. Perach, R. Ronen, and S. Kvatinsky, “Efficient error-correcting-code mechanism for high-throughput memristive processing-in-memory,” in *58th ACM/IEEE Design Automation Conference, DAC 2021, San Francisco, CA, USA, December 5-9, 2021*. IEEE, 2021, pp. 199–204.
  - [25] M. Liu, L. Xia, Y. Wang, and K. Chakrabarty, “Fault tolerance for rram-based matrix operations,” in *IEEE International Test Conference, ITC 2018, Phoenix, AZ, USA, October 29 - Nov. 1, 2018*. IEEE, 2018, pp. 1–10.
  - [26] J. Wang, Q. Xu, B. Yuan, S. Chen, B. Yu, and F. Wu, “Reliability-driven neural network training for memristive crossbar-based neuromorphic computing systems,” in *IEEE International Symposium on Circuits and Systems, ISCAS 2020, Sevilla, Spain, October 10-21, 2020*. IEEE, 2020, pp. 1–4.
  - [27] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
  - [28] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Commun. ACM*, vol. 60, no. 6, pp. 84–90, 2017.
  - [29] G. Yuan, Z. Liao, X. Ma, Y. Cai, Z. Kong, X. Shen, J. Fu, Z. Li, C. Zhang, H. Peng, N. Liu, A. Ren, J. Wang, and Y. Wang, “Improving DNN fault tolerance using weight pruning and differential crossbar mapping for rram-based edge AI,” in *22nd International Symposium on Quality Electronic Design, ISQED 2021, Santa Clara, CA, USA, April 7-9, 2021*. IEEE, 2021, pp. 135–141.
  - [30] E. Ozen and A. Orailoglu, “Architecting decentralization and customizability in DNN accelerators for hardware defect adaptation,” *IEEE Trans. Comput. Aided Des. Integr. Circuits Syst.*, vol. 41, no. 11, pp. 3934–3945, 2022.
  - [31] T. Spyrou, S. A. El-Sayed, E. Afacan, L. A. Camuñas-Mesa, B. Linares-Barranco, and H. Stratigopoulos, “Neuron fault tolerance in spiking neural networks,” in *Design, Automation & Test in Europe Conference & Exhibition, DATE 2021, Grenoble, France, February 1-5, 2021*. IEEE, 2021, pp. 743–748.