

Accelerator Architecture for Secure and Energy Efficient
Machine Learning

Dissertation

Presented in Partial Fulfillment of the Requirements for the Degree
Doctor of Philosophy in the Graduate School of The Ohio State
University

By

Mohammad Hossein Samavatian,

Graduate Program in Department of Computer Science and Engineering

The Ohio State University

2022

Dissertation Committee:

Radu Teoderescu, Advisor

Wei-Lun Chao

Yang Wang

© Copyright by

Mohammad Hossein Samavatian

2022

Abstract

ML applications are driving the next computing revolution. In this context both performance and security are crucial. We propose hardware/software co-design solutions for addressing both. First, we propose RNNFast, an accelerator for Recurrent Neural Networks (RNNs). RNNs are particularly well suited for machine learning problems in which context is important, such as language translation. RNNFast leverages an emerging class of non-volatile memory called domain-wall memory (DWM). We show that DWM is very well suited for RNN acceleration due to its very high density and low read/write energy. RNNFast is very efficient and highly scalable, with a flexible mapping of logical neurons to RNN hardware blocks. The accelerator is designed to minimize data movement by closely interleaving DWM storage and computation. We compare our design with a state-of-the-art GPGPU and find $21.8\times$ higher performance with $70\times$ lower energy.

Second, we brought ML security into ML accelerator design for more efficiency and robustness. Deep Neural Networks (DNNs) are employed in an increasing number of applications, some of which are safety-critical. Unfortunately, DNNs are known to be vulnerable to so-called adversarial attacks. In general, the proposed defenses have high overhead, some require attack-specific re-training of the model or careful tuning to adapt to different attacks. We show that these approaches, while successful

for a range of inputs, are insufficient to address stronger, high-confidence adversarial attacks. To address this, we propose HASI and DNNSHIELD, two hardware-accelerated defenses that adapt the strength of the response to the confidence of the adversarial input. Both techniques rely on approximation or random noise deliberately introduced into the model. HASI uses direct noise injection into the model at inference. DNNSHIELD uses approximation that relies on dynamic and random sparsification of the DNN model to achieve inference approximation efficiently and with fine-grain control over the approximation error. Both techniques use the output distribution characteristics of noisy/sparsified inference compared to a baseline output to detect adversarial inputs. We show an adversarial detection rate of 86% when applied to VGG16 and 88% when applied to ResNet50, which exceeds the detection rate of the state of the art approaches, with a much lower overhead. We demonstrate a software/hardware-accelerated FPGA prototype, which reduces the performance impact of HASI and DNNSHIELD relative to software-only CPU and GPU implementations.

3.10.3 Performance, Area and Power Overheads	108
3.10.4 Sensitivity Studies	110
3.11 Conclusion	115
Bibliography	116

List of Tables

Table	Page
2.1 Summary of racetrack memory parameters.	41
2.2 RNNFast design parameters with associated power and area overheads.	42
2.3 Summary of the benchmarks evaluated.	44
2.4 Energy and run time for FPGA-based RNNs.	52
3.1 Attack parameters for multiple variants of CW and EAD attacks. Original attack success rate, confidence, and distortion. Dataset from ImageNet, on VGG16 and ResNet50.	103
3.2 Detection rates for 4 defenses: SAP, FS, Approximate-MUL, HASI and DNNSHIELD. Dataset from ImageNet, on VGG16 and ResNet50. . .	105
3.3 HASI/DNNSHIELD aware adaptive attacks	108
3.4 FPGA resources and power overhead of DNNSHIELD over baseline CHaiDNN accelerator.	111

3.21 Adversarial detection and benign FPR with different sparsification approaches.	112
3.22 Adversarial attack success rate for multiple attacks as a function of the number of noisy runs in DNNSHIELD.	112
3.23 Number of noisy runs required by DNNSHIELD to make a classification decision for multiple adversarial data sets and benign inputs.	113
3.24 Benign and adversarial FPR for different threshold values, left VGG16 and right ResNet50.	114

Chapter 1: Introduction

Deep learning is transforming the way we approach everyday computing. Deep neural networks (DNNs) are rapidly becoming indispensable tools for solving an increasingly diverse set of complex problems, including computer vision [70], natural language processing [28], machine translation [13], and many others. In this context both performance and security are crucial. Different applications in deep learning have different requirements that may not necessarily be aligned. Some applications need more memory or compute resources depending on the underlying neural network while in other applications privacy and security concerns have higher priority. However, compute efficiency is not necessarily unrelated to security and privacy. For instance, techniques designed to improve energy efficiency can be used indirectly for ML privacy purposes. Mobile devices, while energy constrained, are sufficiently powerful to allow complex ML computation to be performed locally, avoiding the transfer of potentially sensitive user data to the cloud. This does however require energy efficient ML solutions. This is especially important for memory-intensive ML models such as recurrent neural networks.

Applications like speech recognition empower today’s digital assistants, business intelligence applications fueled by the analysis of social media postings, etc. For these applications, processing information in a way that preserves the correct context

Chapter 2: RNNFast: An Accelerator for Recurrent Neural Networks Using Domain Wall Memory

Recurrent Neural Networks (RNNs) are a powerful class of networks designed to consider context by retaining and using information from previously processed inputs. RNNs have the ability to learn sequences and can be applied to any problems that require context that needs to be remembered. RNNs are used across a wide range of applications that include speech recognition for digital assistants such as Siri and Google Now, sentiment analysis for classifying social media postings, and language translation.

However, RNN workloads are data-intensive because they store a partial history of the output sequence and perform computations on that history along with the current input. As a result, RNNs require both vast amounts of storage and increased processing power. For example, the RNN neuron requires $8\times$ the number of weights and multiply-accumulate (MAC) operations of a typical FC neuron. RNN networks are also generally quite large. For instance, Amodei et al. [7] developed a network for performing speech recognition that utilized seven recurrent layers and a total of 35 million parameters. At this scale, RNNs with large input sets are susceptible to memory bottlenecks when running on existing accelerators such as GPUs [49] or FPGAs [49, 76, 38, 91, 11, 135, 77, 136]. In addition, the fundamentally different

overshift is on the order of 10^{-5} [149], which is quite high. However, the probability of multibit overshift is about 10^{-21} , which is negligible. As a result, RNNFast implements mitigation for single-bit overshift errors.

through a crossbar. These four gate modules perform LSTM-RNN inference, and transport results to LSTM Functional Logic to perform the remaining computation (element-wise multiplication and addition of gate vectors, activations, etc.). Then, the final results are loaded to output buffer groups through a crossbar. The current state of the LSTM cell is stored in an on-chip buffer, called Cell Buffer.

Inside each gate module, gate vector is calculated in a tiling scheme. Tiled input vectors and the corresponding parameters are transferred into the LSTM gate module in parallel to perform inference. Inside each LSTM gate module, all multiplications between input elements and parameters are performed in parallel. The results are then summed up through an addition tree to minimize latency. The whole architecture is also pipelined to further improve throughput. The outputs are fed into activation nodes to generate the final output vectors of each gate.

Domain-Wall Memory is an increasingly popular candidate for replacing conventional memories such as Flash, DRAM and SRAM, and there are prior works utilizing DWM in reconfigurable computing and machine learning architectures designs [139, 140, 147, 60, 26, 153]. Zhao et al. [153] employed racetrack memory for reconfigurable computing to achieve high density and low energy compared with SRAM. Chung et al. [26] proposed a DWM dot product engine using a DWM-based analog design, which requires ADCs. Yu et al. [147] designed data intensive machine learning image-processing into in-memory DWM. The high storage density offered by racetrack memory makes it a promising candidate for the data-intensive machine learning applications.

K80. RNNFast is up to $260\times$ faster than the newer NVIDIA P100 for workloads of similar size.

2.7 Conclusion

The unprecedented growth of available data is accelerating the adoption of deep learning across a wide range of applications including speech recognition, machine translation, and language modeling. In this study, we present RNNFast, a novel accelerator designed for recurrent neural networks. Our design demonstrates that using domain wall memory is not only feasible, but also very efficient. We compare RNNFast with a state-of-the-art P100 NVIDIA GPU and find $21.8\times$ better performance with $70\times$ lower energy.

Chapter 3: Accelerator Architecture for ML Security

Convolutional neural networks have demonstrated high accuracy on various tasks in recent years. However they are extremely vulnerable to adversarial examples. For example, imperceptible perturbations added to clean images can cause convolutional network to fail. Figure 3.1 shows two examples of adversarial images generated using the state-of-the-art $CW-L_2$ attack [18]. The leftmost images are benign, unmodified samples. They are correctly classified by a DNN model such as VGG16 with 99% and 87% confidence, respectively. The middle and rightmost pairs of images represent the output of two versions of the $CW-L_2$ attack, each resulting in misclassification. Note that all adversarial images are virtually indistinguishable from the original to the casual observer, even though the confidence of the classifier in all cases is very high.

Several defenses have been proposed to address adversarial attacks [86, 102, 146, 34, 15, 87]. Most rely on purely software implementations, with high overheads, limiting their utility to real-world applications. A recent line of research has explored hardware-assisted approximate computing to introduce controlled errors into the inference process, either through model quantization [41, 101] or approximate computation [50]. This inference approximation disrupts the effect of the adversarial

detection rate of the state of the art approaches. We also show that DNNShield is robust against attacks that are aware of our defense and attempt to circumvent it.

The accelerator design builds explicit support for dynamic and random model sparsification. The DNNShield accelerator is optimized for efficiently executing sparsified models in which the sparsification rate changes as a function of the input – which is more challenging compared to models for which weight sparsity is fixed. We show that the DNNShield accelerator reduces the performance impact of approximate inference-based adversarial detection to $1.53 \times -2 \times$ relative to the unprotected baseline, compared to $15 \times -25 \times$ overhead for a software-only GPU implementation.

The rest of this chapter is organized as follows: Section 3.1 provides background information. Section 3.2 discusses related work. Section 3.3 explains the threat model. Section 3.4 focuses on high confident adversarial examples and how prior approximate methods fail to detect such anomalies. Section 3.5 presents the fundamental of the detection mechanism used in our work. Section 3.6 and 3.7 present the details of the HASI and DNNShield designs. Section 3.8 provide the information of the FPGA implementation of both designs. Finally section 3.9 and 3.10 shows the evaluation results and Section 3.11 concludes the chapter.

correct classification in the presence of adversarial inputs; or (b) aim to provide detection, rejecting suspicious inputs. We will discuss some of the important works in section 3.2.

3.3 Threat Model

We assume in this work that the adversary has complete access to the network, including the output prediction and logits, with full knowledge of the architecture and parameters, and is able to use this in a white-box manner. We focus mainly on state-of-the-art optimization-based attacks—CW and EAD—since it has been demonstrated that all earlier attacks can be overcome utilizing other methods, such as adversarial training [45] or defensive distillation [102], which could be used in combination with our approach. Additionally, we verify our evaluation includes high confidence adversarial examples, as some previously proposed defenses were later shown to perform poorly under a more holistic treatment which included these [85].

3.4 Motivation

Strong attacks such as CW and EAD can be tuned to produce a class of adversarial inputs that present a significant challenge to approximation-based defenses. Prior work has shown that adversarial inputs can be constructed to induce misclassification with very high confidence [119, 16, 18]. In other words, the victim model assigns a very high probability to the adversarial input belonging to the wrong class. These so-called “high confidence” adversarials can be constructed while minimizing the distortion to original input.

3.4.1 High-Confidence Attack Variants

Figure 3.2 shows an example of multiple adversarial samples for a benign image with different levels of classification confidence and the corresponding distortion. To measure classification confidence we used the Z-score (the number of standard deviations by which the value of a raw score is above or below the mean value) of the maximum logit value, which corresponds to the class with the highest confidence. Adv_1 is a low distortion adversarial of the benign with low classification confidence of 4.18. Adv_2 is a high confidence example of the same input with very high classification confidence of 12. While distortion is also higher, it is still imperceptible to the naked eye. We will show that existing defenses are ineffective against this type of adversarial. Increasing the confidence beyond 12 increased distortion significantly, as Adv_3 shows.

- [119] Yash Sharma and Pin-Yu Chen. Bypassing feature squeezing by increasing adversary strength. *arXiv preprint arXiv:1803.09868*, 2018.
- [120] Yongming Shen, Michael Ferdman, and Peter Milder. Maximizing cnn accelerator efficiency through resource partitioning. In *Proceedings of the 44th Annual International Symposium on Computer Architecture, ISCA '17*, pages 535–547, New York, NY, USA, 2017. ACM.
- [121] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2014.
- [122] Clinton W. Smullen, Anurag Nigam, Sudhanva Gurumurthi, and Mircea R. Stan. The stesims stt-ram simulation and modeling system. In *iccad*, pages 318–325. IEEE Press, 2011.
- [123] Z. Sun, X. Bi, W. Wu, S. Yoo, and H. (. Li. Array organization and data management exploration in racetrack memory. *IEEE Transactions on Computers*, 65(4):1041–1054, April 2016.
- [124] Zhanrui Sun, Yongxin Zhu, Yu Zheng, Hao Wu, Zihao Cao, Peng Xiong, Junjie Hou, Tian Huang, and Zhiqiang Que. Fpga acceleration of lstm based on data for test flight. In *2018 IEEE International Conference on Smart Cloud (SmartCloud)*, pages 1–6. IEEE, 2018.
- [125] Zhenyu Sun, Wenqing Wu, and Hai (Helen) Li. Cross-layer racetrack memory design for ultra high density and low power consumption. In *Proceedings of the 50th Annual Design Automation Conference, DAC '13*, pages 53:1–53:6, New York, NY, USA, 2013. ACM.
- [126] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations 2014*, 2014.
- [127] MT Tommiska. Efficient digital implementation of the sigmoid function for reprogrammable logic. *IEE Proceedings-Computers and Digital Techniques*, 150(6):403–411, 2003.
- [128] Antonio Toral and Andy Way. What level of quality can neural machine translation attain on literary text? In *Translation Quality Assessment*, pages 263–287. Springer, 2018.
- [129] Florian Tramèr, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. *arXiv preprint arXiv:2002.08347*, 2020.

- [130] F. Vaverka, V. Mrazek, Z. Vasicek, L. Sekanina, M. A. Hanif, and M. Shafique. Tfapprox: Towards a fast emulation of dnn approximate hardware accelerators on gpu. In *2020 Design, Automation and Test in Europe Conference (DATE)*, page 4, 2020.
- [131] Swagath Venkataramani, Ashish Ranjan, Subarno Banerjee, Dipankar Das, Sasikanth Avancha, Ashok Jagannathan, Ajaya Durg, Dheemanth Nagaraj, Bharat Kaul, Pradeep Dubey, and Anand Raghunathan. Scaleddeep: A scalable compute architecture for learning and evaluating deep networks. In *Proceedings of the 44th Annual International Symposium on Computer Architecture, ISCA '17*, pages 13–26, New York, NY, USA, 2017. ACM.
- [132] R. Venkatesan, M. Sharad, K. Roy, and A. Raghunathan. Dwm-tapestri - an energy efficient all-spin cache using domain wall shift based writes. In *2013 Design, Automation Test in Europe Conference Exhibition (DATE)*, pages 1825–1830, March 2013.
- [133] Rangharajan Venkatesan, Vivek J. Kozhikkottu, Mrigank Sharad, Charles Augustine, Arijit Raychowdhury, Kaushik Roy, and Anand Raghunathan. Cache design with domain wall memory. *IEEE Trans. Comput.*, 65(4):1010–1024, April 2016.
- [134] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: Lessons learned from the 2015 MSCOCO image captioning challenge. *CoRR*, abs/1609.06647, 2016.
- [135] M. Wang, Z. Wang, J. Lu, J. Lin, and Z. Wang. E-lstm: An efficient hardware architecture for long short-term memory. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, pages 1–1, 2019.
- [136] Shuo Wang, Zhe Li, Caiwen Ding, Bo Yuan, Qinru Qiu, Yanzhi Wang, and Yun Liang. C-lstm: Enabling efficient lstm using structured compression techniques on fpgas. In *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, FPGA '18*, pages 11–20, New York, NY, USA, 2018. ACM.
- [137] X. Wang, J. Yu, C. Augustine, R. Iyer, and R. Das. Bit prudent in-cache acceleration of deep convolutional neural networks. In *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 81–93, Feb 2019.
- [138] Xingbin Wang, Rui Hou, Boyan Zhao, Fengkai Yuan, Jun Zhang, Dan Meng, and Xuehai Qian. Dnnguard: An elastic heterogeneous dnn accelerator architecture against adversarial attacks. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages*

and Operating Systems, ASPLOS '20, page 19–34, New York, NY, USA, 2020. Association for Computing Machinery.

- [139] Yuhao Wang, Hao Yu, Leibin Ni, Guang-Bin Huang, Mei Yan, Chuliang Weng, Wei Yang, and Junfeng Zhao. An energy-efficient nonvolatile in-memory computing architecture for extreme learning machine by domain-wall nanowire devices. *IEEE Transactions on Nanotechnology*, 14(6):998–1012, 2015.
- [140] Yuhao Wang, Hao Yu, Dennis Sylvester, and Pingfan Kong. Energy efficient in-memory AES encryption based on nonvolatile domain-wall nanowire. In *Design, Automation & Test in Europe Conference & Exhibition, DATE 2014, Dresden, Germany, March 24-28, 2014*, pages 1–4, 2014.
- [141] Yulong Wang, Hang Su, Bo Zhang, and Xiaolin Hu. Interpret neural networks by identifying critical data routing paths. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8906–8914, 2018.
- [142] Zhisheng Wang, Jun Lin, and Zhongfeng Wang. Accelerating recurrent neural networks: A memory-efficient approach. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 25(10):2763–2775, 2017.
- [143] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization. In *International Conference on Learning Representations (ICLR)*, 2018.
- [144] Xilinx. CHaiDNN. <https://github.com/Xilinx/CHaiDNN>.
- [145] Cong Xu, Dimin Niu, Xiaochun Zhu, Seung H. Kang, Matt Nowak, and Yuan Xie. Device-architecture co-optimization of STT-RAM based memory for low power embedded systems. In *iccad*, pages 463–470. IEEE Press, 2011.
- [146] Weilin Xu, David Evans, and Yanjun Qi. Feature squeezing: Detecting adversarial examples in deep neural networks. In *Proceedings 2018 Network and Distributed System Security Symposium*, 2018.
- [147] Hao Yu, Yuhao Wang, Shuai Chen, Wei Fei, Chuliang Weng, Junfeng Zhao, and Zhulin Wei. Energy efficient in-memory machine learning for data intensive image-processing by non-volatile domain-wall memory. In *19th Asia and South Pacific Design Automation Conference, ASP-DAC 2014, Singapore, January 20-23, 2014*, pages 191–196, 2014.
- [148] Chao Zhang, Guangyu Sun, Weiqi Zhang, Fan Mi, Hai Li, and W. Zhao. Quantitative modeling of racetrack memory, a tradeoff among area, performance, and power. In *The 20th Asia and South Pacific Design Automation Conference*, pages 100–105, Jan 2015.

- [149] Chao Zhang, Guangyu Sun, Xian Zhang, Weiqi Zhang, Weisheng Zhao, Tao Wang, Yun Liang, Yongpan Liu, Yu Wang, and Jiwu Shu. Hi-fi playback: Tolerating position errors in shift operations of racetrack memory. In *Proceedings of the 42Nd Annual International Symposium on Computer Architecture, ISCA '15*, pages 694–706, New York, NY, USA, 2015. ACM.
- [150] Shijin Zhang, Zidong Du, Lei Zhang, Huiying Lan, Shaoli Liu, Ling Li, Qi Guo, Tianshi Chen, and Yunji Chen. Cambricon-x: An accelerator for sparse neural networks. In *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 1–12. IEEE, 2016.
- [151] Y. Zhang, C. Zhang, J. Nan, Z. Zhang, X. Zhang, J. O. Klein, D. Ravelosona, G. Sun, and W. Zhao. Perspectives of racetrack memory for large-capacity on-chip memory: From device to system. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 63(5):629–638, May 2016.
- [152] Yiwei Zhang, Chao Wang, Lei Gong, Yuntao Lu, Fan Sun, Chongchong Xu, Xi Li, and Xuehai Zhou. A power-efficient accelerator based on fpgas for lstm network. In *2017 IEEE International Conference on Cluster Computing (CLUSTER)*, pages 629–630. IEEE, 2017.
- [153] Weisheng Zhao, Nesrine Ben Romdhane, Yue Zhang, Jacques-Olivier Klein, and Define Ravelosona. Racetrack memory based reconfigurable computing. In *Faible Tension Faible Consommation (FTFC), 2013 IEEE*, pages 1–4. IEEE, 2013.