

Robust and Energy-Efficient Deep Learning Systems

Muhammad Abdullah Hanif¹ (Ph.D. Candidate), Muhammad Shafique² (Advisor)

¹Technische Universität Wien (TU Wien), Vienna, Austria

²Division of Engineering, New York University Abu Dhabi (NYUAD), Abu Dhabi, UAE

Problems and Motivation

Applications



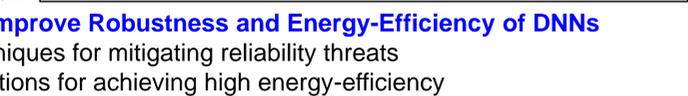
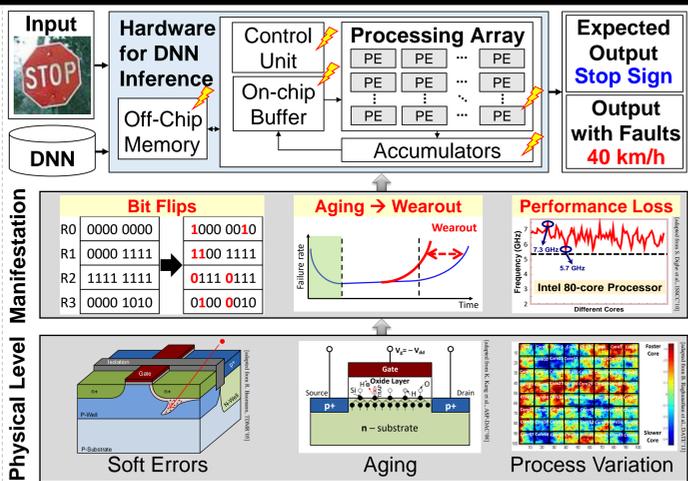
Image Classification, Natural Object Detection & Localization, etc. Language Processing

Compute Cost [OPs]

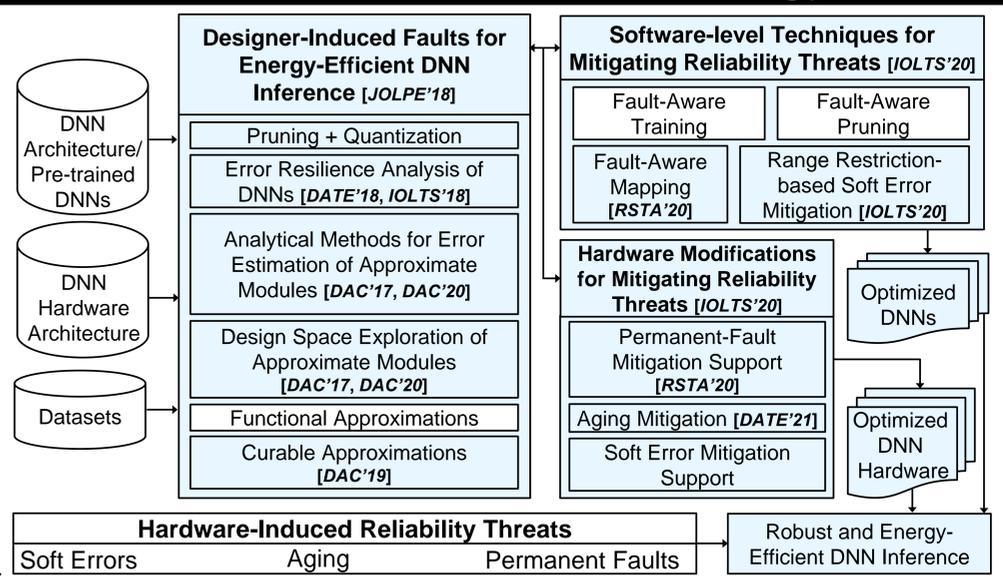
- AlexNet → 1.5B
- VGGNet → 19.6B
- Inception → 2B
- ResNet-152 → 11B

Design Methods to Improve Robustness and Energy-Efficiency of DNNs

- Cost-effective techniques for mitigating reliability threats
- HW/SW approximations for achieving high energy-efficiency

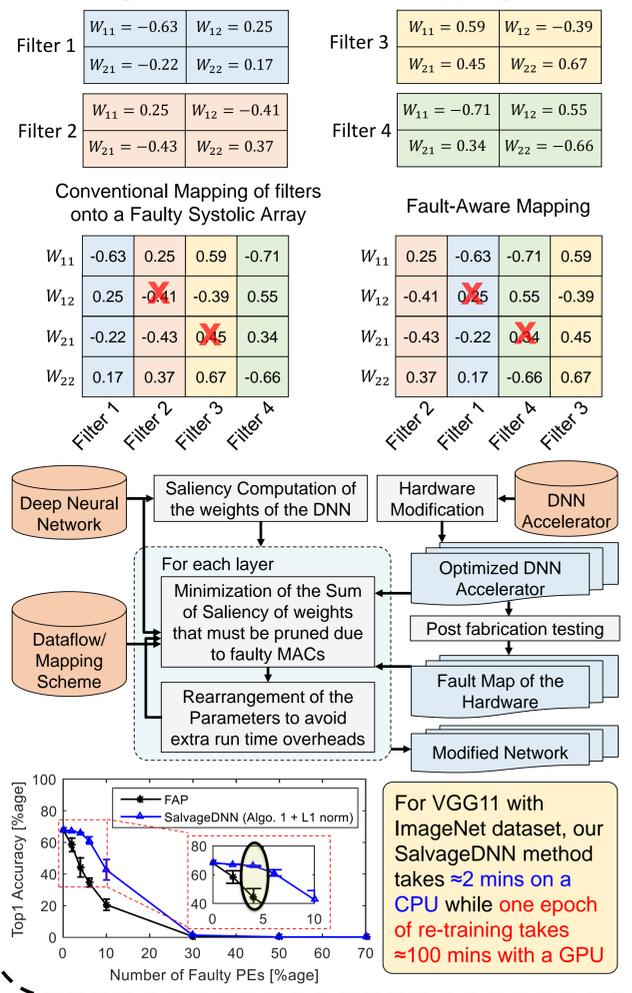


Overview of Our Methodology

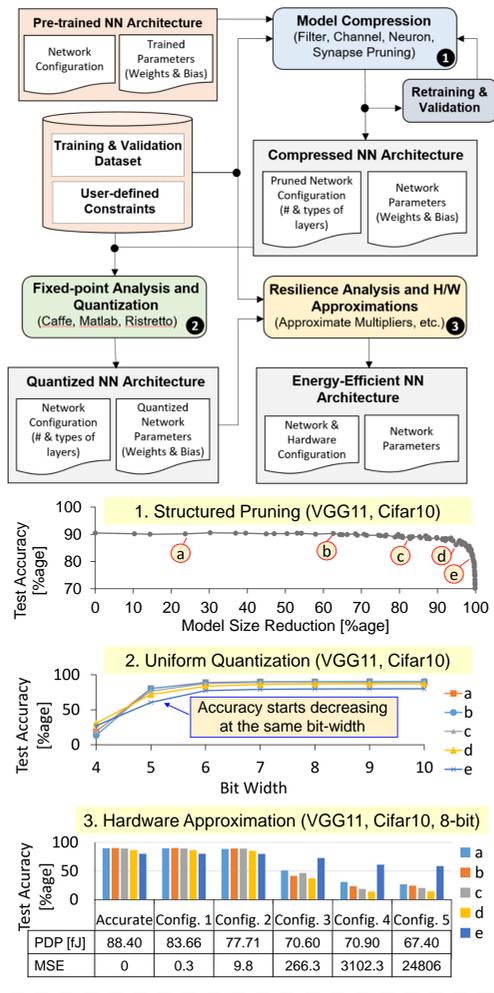


Techniques for Robust and Energy-Efficient Deep Neural Network Inference

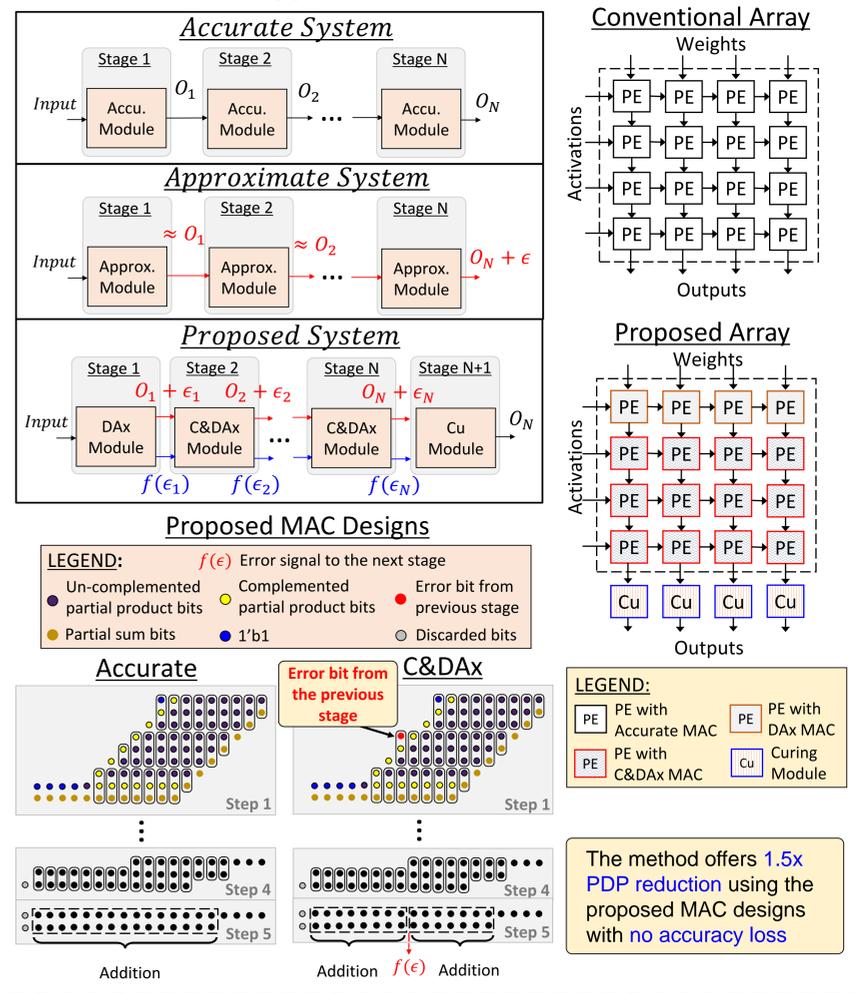
SalvageDNN: Fault-Aware Mapping [RSTA'20]



HW/SW Approximations [JOLPE'18]



CANN: Curable Approximations [DAC'19]



Selected Publications

[DAC'17] M. A. Hanif, R. Hafiz, O. Hasan, M. Shafique, "QuAd: Design and analysis of quality-area optimal low-latency approximate adders", Design Automation Conference (DAC), pp. 1–6, 2017. **Received a HiPEAC Paper Award.**

[DATE'18] M. A. Hanif, R. Hafiz, M. Shafique, "Error resilience analysis for systematically employing approximate computing in convolutional neural networks", DATE, pp. 913–916, 2018.

[DAC'19] M. A. Hanif, F. Khalid, M. Shafique, "CANN: Curable approximations for high-performance deep neural network accelerators", DAC, pp. 1–6, 2019. **Received a HiPEAC Paper Award.**

[DAC'20] M. A. Hanif, R. Hafiz, O. Hasan, M. Shafique, "PEMACx: A probabilistic error analysis methodology for adders with cascaded approximate units", DAC, pp. 1–6, 2020. **Received a HiPEAC Paper Award.**

[DATE'21] M. A. Hanif, M. Shafique, "DNN-Life: An energy-efficient aging mitigation framework for improving the lifetime of on-chip weight memories in deep neural network hardware architectures", DATE, pp. 1-6, 2021.

[JOLPE'18] M. A. Hanif, A. Marchisio, T. Arif, R. Hafiz, S. Rehman, M. Shafique, "X-DNNs: Systematic cross-layer approximations for energy-efficient deep neural networks", Journal of Low Power Electronics (JOLPE), vol. 14, no. 4, pp. 520–534, 2018.

[IOLTS'18] M. A. Hanif, F. Khalid, R. V. W. Putra, S. Rehman, M. Shafique, "Robust machine learning systems: Reliability and security for deep neural networks", IEEE International Symposium on On-Line Testing and Robust System Design (IOLTS), pp. 257–260, 2018.

[IOLTS'20] M. A. Hanif, M. Shafique, "Dependable deep learning: Towards cost-efficient resilience of deep neural network accelerators against soft errors and permanent faults", IOLTS, pp. 1–4, 2020.

[RSTA'20] M. A. Hanif, M. Shafique, "SalvageDNN: Salvaging deep neural network accelerators with permanent faults through saliency-driven fault-aware mapping", Philosophical Transactions of the Royal Society A (RSTA), vol. 378, no. 2164, pp. 20190164, 2020.