

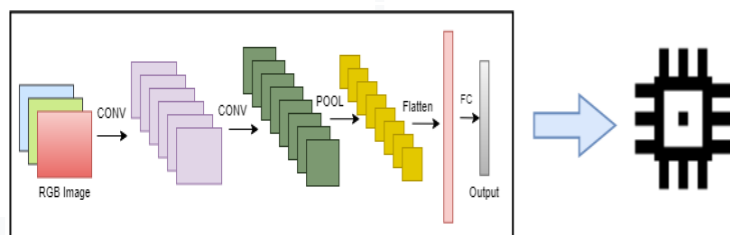
# Adversarial Robust Model Compression Using In-train Pruning

Manoj-Rohit Vemparala<sup>1</sup>, Nael Fasfous<sup>2</sup>, Alexander Frickenstein<sup>1</sup>, Sreetama Sarkar<sup>1</sup>, Qi Zhao<sup>3</sup>, Sabine Kuhn<sup>1</sup>, Lukas Frickenstein<sup>1</sup>, Anmol Singh<sup>1</sup>, Christian Unger<sup>1</sup>, Naveen-Shankar Nagaraja<sup>1</sup>, Christian Wressnegger<sup>3</sup>, Walter Stechele<sup>2</sup>  
<sup>1</sup> BMW Autonomous Driving, <sup>2</sup> Technical University of Munich, <sup>3</sup> Karlsruhe Institute of Technology

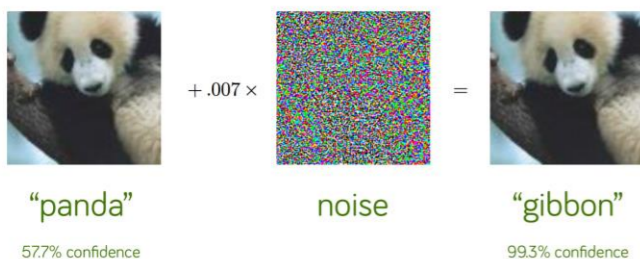
## Motivation

**Goal:** Secure deployment of CNNs on edge devices

**Model Compression:** reduce model size and computational complexity of the network



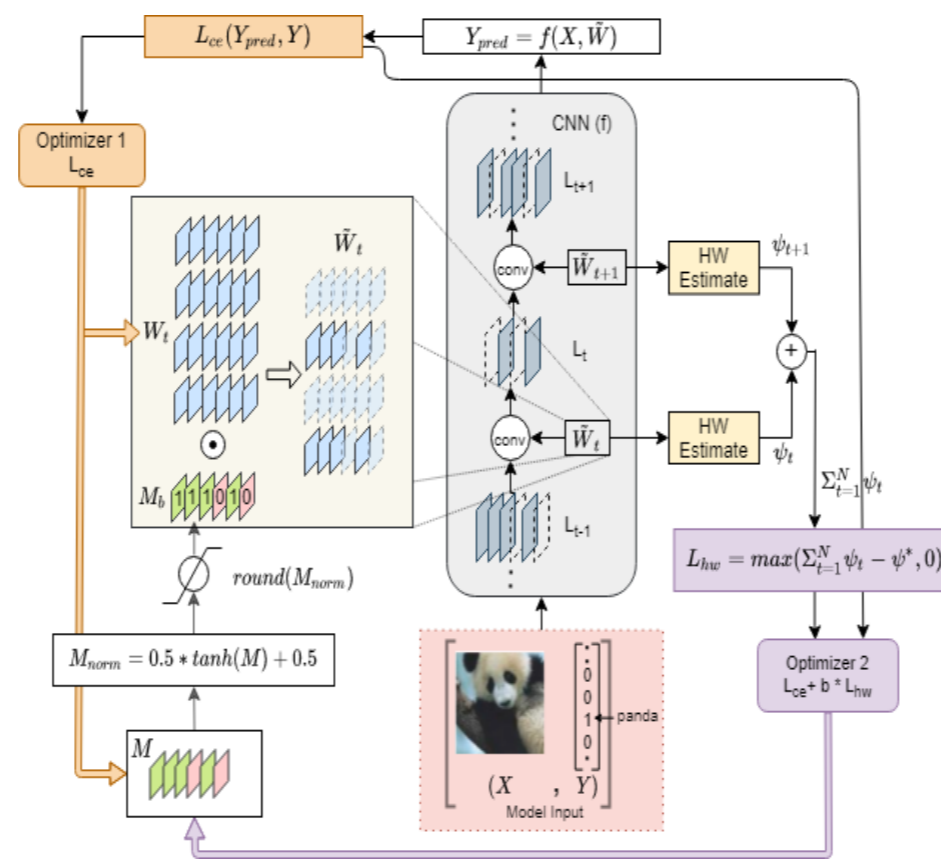
**Adversarial Robustness:** correctly classify images generated using adversarial perturbations



Goodfellow, Ian J. et al. "Explaining and Harnessing Adversarial Examples." CoRR (2015)

## Our Solution: In-train Pruning

- Introduce trainable prune masks which are trained along with network weights
- Training prune masks jointly optimizes cross-entropy loss and hardware loss  $L_{hw}$
- For robust pruning, Fast Adversarial Training is used in place of normal training



## Experimental Results

- This method alleviates the need for *pre-trained model* and *post-train pruning*
- Improves natural accuracy while maintaining same level of adversarial robustness compared to Sota methods

| Method                   | Model Name | Model Size | Natural Acc (%) | PGD Acc (%)  |
|--------------------------|------------|------------|-----------------|--------------|
| Robust ADMM <sup>1</sup> | ResNet18   | 0.04       | 64.52           | 38.01        |
| Ours                     | ResNet20   | 0.04       | <b>70.73</b>    | <b>39.31</b> |
| Robust ADMM <sup>1</sup> | ResNet18   | 0.17       | 73.35           | 43.17        |
| Ours                     | ResNet20   | 0.16       | <b>79.67</b>    | <b>43.22</b> |

<sup>1</sup>Ye, Shaokai et al. "Adversarial Robustness vs. Model Compression, or Both?" 2019 IEEE/CVF International Conference on Computer Vision (ICCV) (2019)