

Model Compression with Adversarial Robustness: A Unified Optimization Framework

Shupeng Gui_{1,†}, Haotao Wang_{2,†}, Haichuan Yang₁, Chen Yu₁, Zhangyang Wang₂, and Ji Liu₃

₁Department of Computer Science, University of Rochester

₂Department of Computer Science and Engineering, Texas A&M University

₃Ytech Seattle AI lab, FeDA lab, AI platform, Kwai Inc

† Equal contribution

NeurIPS 2019, Vancouver CANADA



UNIVERSITY of
ROCHESTER

VITA

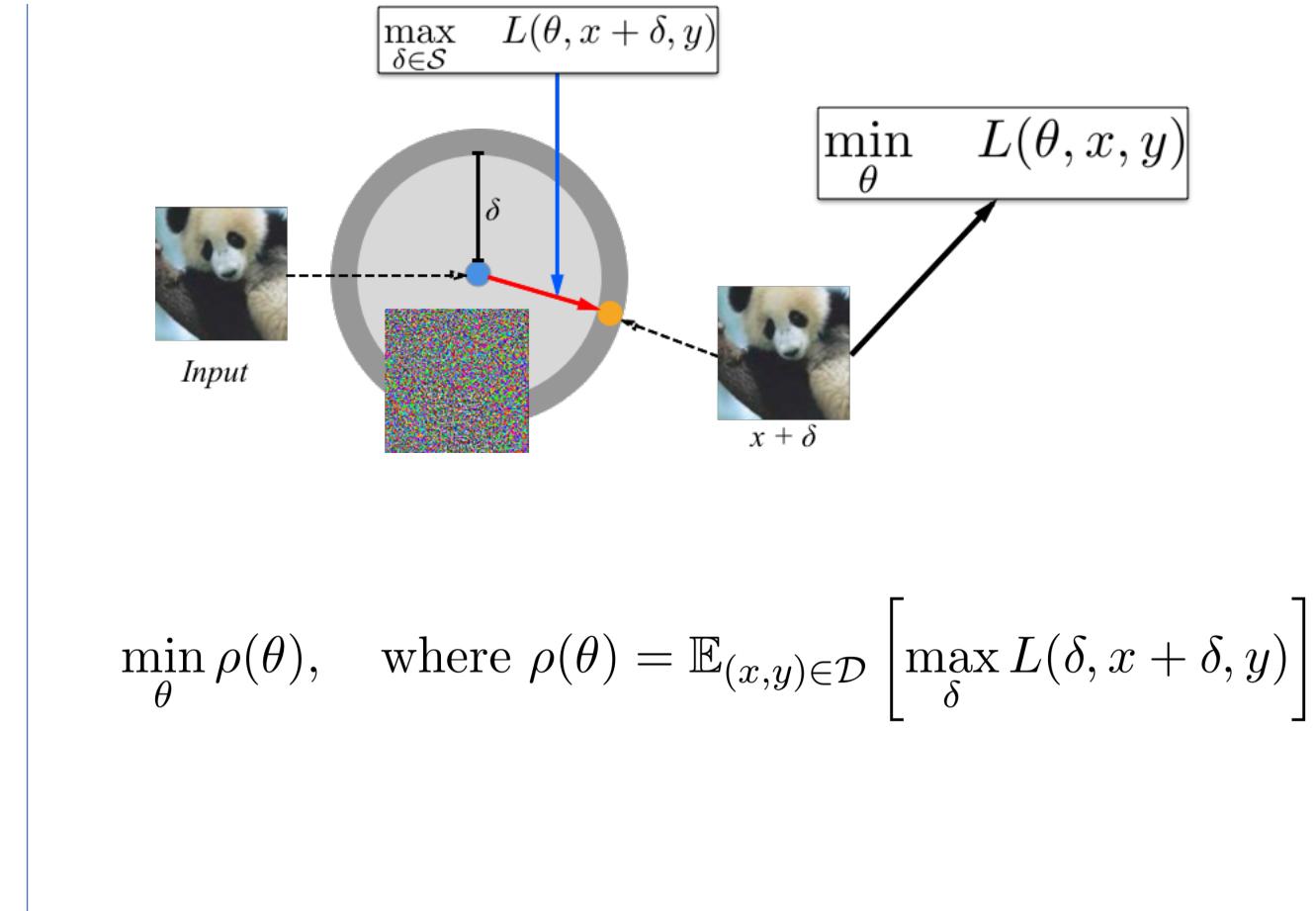
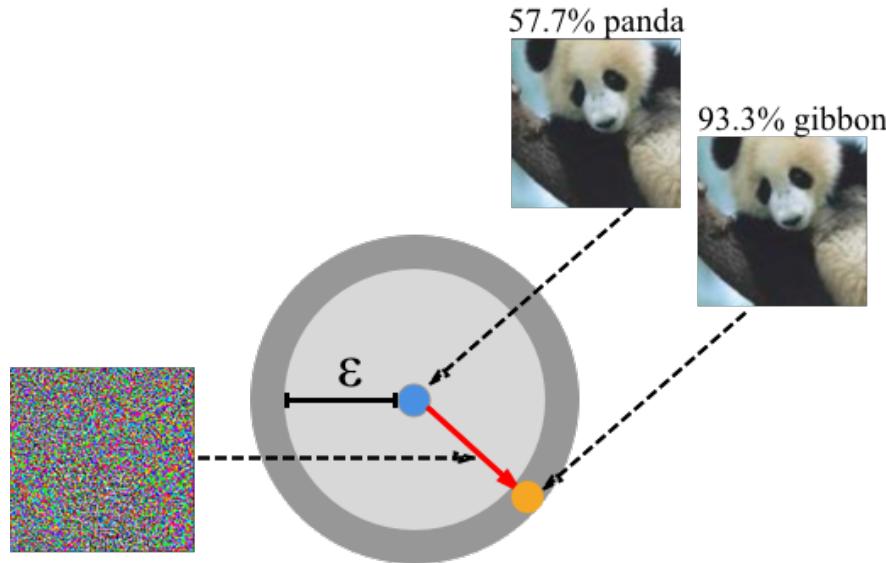
ATM | TEXAS A&M
UNIVERSITY®

OG

Highlights

- Compressing models without hurting their robustness to adversarial attacks, in addition to maintaining accuracy
- Adversarially Trained Model Compression (ATMC) framework
- Integrating pruning, factorization, and quantization into constraints
- An extensive group of experiments demonstrate that ATMC obtains remarkably more favorable trade-off among model size, accuracy and robustness, over currently available alternatives in various settings.

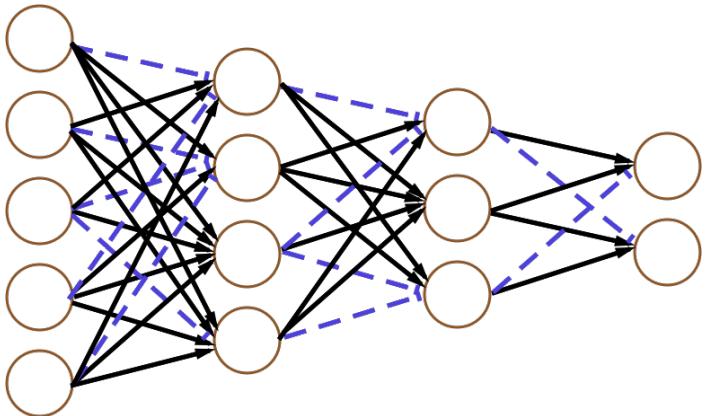
Robustness for CNN



Goodfellow et al, “Explaining and Harnessing Adversarial Examples”, ICLR 2015.

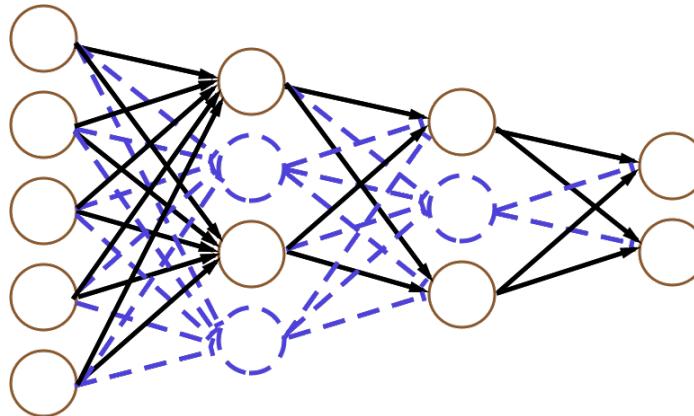
Madry, Aleksander, et al. “Towards deep learning models resistant to adversarial attacks.”, ICLR 2018

Compression for CNN



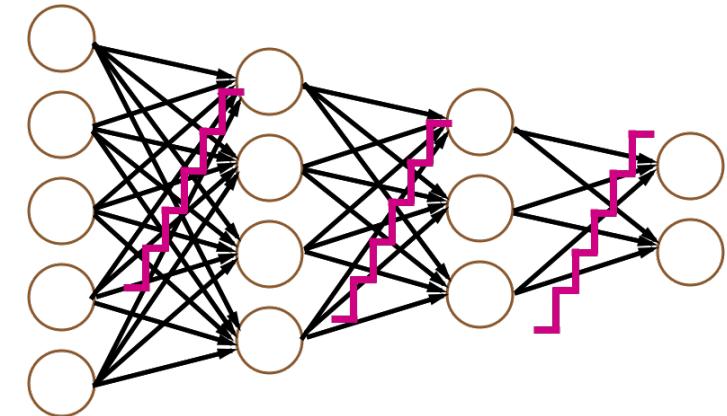
Weight Pruning

[Han et al., NIPS 15]



Node Pruning

[Li et al., ICLR 17]
[He et al., ICCV 17]
[Porikli et al., ECCV 16]



Quantization

[Han et al., ICLR 16]

Compression v.s. Robustness

- Highly non-straightforward and contextually varying w.r.t different means of compression
- An appropriately higher CNN model sparsity led to better robustness, whereas over-sparsification could in turn cause more fragility [1]
- [2] showed that sparse algorithms are not stable:
 - If an algorithm promotes sparsity, then its sensitivity to small perturbations for the input data remains bounded away from zero
 - Quantization seems to reduce the Minimum Description Length and might potentially make the algorithm more robust.
- [3] argued that the tradeoff between robustness and accuracy may be inevitable for the classification task.

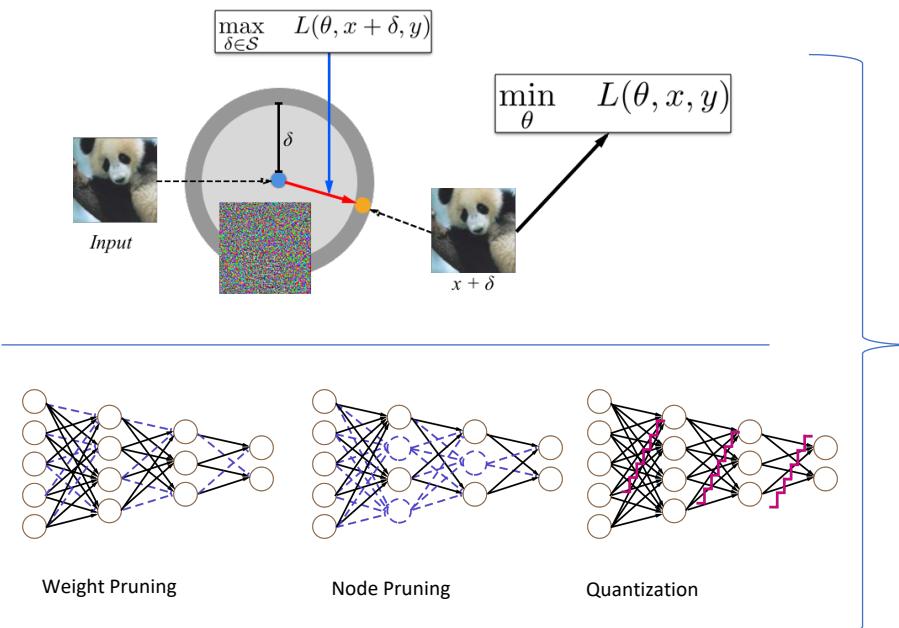
[1] Guo et al, “Sparse dnns with improved adversarial robustness”, NeurIPS 2018

[2] Xu et al, “Sparse algorithms are not stable: A no-free-lunch theorem”. IEEE t-PAMI 2012

[3] Tsipras et al, “Robustness may be at odds with accuracy”. Stat, 1050:11, 2018

Adversarially Trained Model Compression (Overall)

- A unified optimization framework for both compression and robustness on CNN



$$\begin{aligned} & \min_{\theta} \sum_{(x,y) \in \mathcal{Z}} f^{\text{adv}}(\theta; x, y) \\ \text{s.t. } & \underbrace{\sum_{i=1}^l \|\boldsymbol{U}^{(i)}\|_0 + \|\boldsymbol{V}^{(i)}\|_0 + \|\boldsymbol{C}^{(i)}\|_0}_{\|\theta\|_0} \leq k, \text{ (sparsity constraint)} \\ & \theta \in \mathcal{Q}_b := \{\theta : |\boldsymbol{U}^{(l)}|_0 \leq 2^b, |\boldsymbol{V}^{(l)}|_0 \leq 2^b, |\boldsymbol{C}^{(l)}|_0 \leq 2^b \forall l \in [L]\}. \text{ (quantization constraint)} \end{aligned} \quad (4)$$

Adversarially Trained Model Compression (Unified Pruning Constraints)

- Sparsification and channel pruning into one constraint

$$\mathbf{W} = \mathbf{U}\mathbf{V} + \mathbf{C}, \quad \|\mathbf{U}\|_0 + \|\mathbf{V}\|_0 + \|\mathbf{C}\|_0 \leq k$$

where $\|\cdot\|_0$ denotes the number of nonzeros of the augment matrix

- The above enforces a novel, compound (including both multiplicative and additive) sparsity structure on \mathbf{W}
- $\mathbf{U}\mathbf{V}$ indicates a natural form of channel pruning
- Such form can be used in the model acceleration with other norm such as L21

Adversarially Trained Model Compression (Non-uniform Quantization)

- Each nonzero element of the DNN parameter can only be chosen from a set of a few values
- These values are not necessarily evenly distributed and need to be optimized
- With respect to the non-uniform quantization strategy, we introduce the quantization constraints

$$|\mathbf{U}^{(l)}|_0 \leq 2^b, |\mathbf{V}^{(l)}|_0 \leq 2^b, |\mathbf{C}^{(l)}|_0 \leq 2^b \quad \forall l \in [L].$$

Adversarially Trained Model Compression (Adversarial Robustness)

- White-box attack setting: an adversary to eavesdrop the optimization and gradients of the learning model
- When an “clean” image x comes to a target model, “perturb” the image with an adversarial perturbation with bounded magnitudes.

$$B_\infty^\Delta(x) := \{x' : \|x' - x\|_\infty \leq \Delta\}$$

- The training objective for the attacker

$$f^{\text{adv}}(\boldsymbol{\theta}; x, y) = \max_{x' \in B_\infty^\Delta(x)} f(\boldsymbol{\theta}; x', y)$$

- Defensive objective for ATMC

$$\min_{\boldsymbol{\theta}} \sum_{(x,y) \in \mathcal{Z}} f^{\text{adv}}(\boldsymbol{\theta}; x, y).$$

Optimization for ATMC

- Apply ADMM to solve the optimization problem
- Update x^{adv} : $x^{adv} \leftarrow \text{Proj}_{\{x': \|x' - x\|_\infty \leq \Delta\}} \{x + \alpha \nabla_x f(\boldsymbol{\theta}; x, y)\}$
- Update $\boldsymbol{\theta}$: $\min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}; x^{adv}, y) + \frac{\rho}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}' + \mathbf{u}\|_F^2 \quad \text{s.t. } \|\boldsymbol{\theta}\|_0 \leq k.$
 - Only related to the sparsity constraint
$$\boldsymbol{\theta} \leftarrow \text{Proj}_{\{\boldsymbol{\theta}'': \|\boldsymbol{\theta}''\|_0 \leq k\}} \left(\boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} \left[f(\boldsymbol{\theta}; x^{adv}, y) + \frac{\rho}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}' + \mathbf{u}\|_F^2 \right] \right)$$
- Update $\boldsymbol{\theta}'$: $\min_{\boldsymbol{\theta}'} \|\boldsymbol{\theta}' - (\boldsymbol{\theta} + \mathbf{u})\|_F^2, \quad \text{s.t. } \boldsymbol{\theta}' \in \mathcal{Q}_b$
 - Only related to quantization constraint
 - For example, essentially solving
$$\min_{\mathbf{U}, \{\mathbf{a}_k\}_{k=1}^{2^b}} \|\mathbf{U} - \bar{\mathbf{U}}\|^2 \quad \text{s.t. } \mathbf{U}_{i,j} \in \{0, \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{2^b}\}$$
- Lloyd's algorithm [1]
$$\mathbf{U}_t'^{(l)} = \text{ZeroKmeans}_{2^b}(\mathbf{U}^{(l)} + \mathbf{u}_{\mathbf{U}^{(l)}})$$

Optimization for ATMC

Algorithm 1 ZeroKmeans_B(\bar{U})

```
1: Input: a set of real numbers  $\bar{U}$ , number  
   of clusters  $B$ .  
2: Output: quantized tensor  $U$ .  
3: Initialize  $a_1, a_2, \dots, a_B$  by randomly  
   picked nonzero elements from  $\bar{U}$ .  
4:  $Q := \{0, a_1, a_2, \dots, a_B\}$   
5: repeat  
6:   for  $i = 0$  to  $|\bar{U}| - 1$  do  
7:      $\delta_i \leftarrow \arg \min_j (\bar{U}_i - Q_j)^2$   
8:   end for  
9:   Fix  $Q_0 = 0$   
10:  for  $j = 1$  to  $B$  do  
11:     $a_j \leftarrow \frac{\sum_i \mathbf{I}(\delta_i=j) \bar{U}_i}{\sum_i \mathbf{I}(\delta_i=j)}$   
12:  end for  
13: until Convergence  
14: for  $i = 0$  to  $|\bar{U}| - 1$  do  
15:    $U_i \leftarrow Q_{\delta_i}$   
16: end for
```

Algorithm 2 ATMC

```
1: Input: dataset  $\mathcal{Z}$ , stepsize sequence  
    $\{\gamma_t > 0\}_{t=0}^{T-1}$ , update steps  $n$  and  $T$ ,  
   hyper-parameter  $\rho, k$ , and  $b, \Delta$   
2: Output: model  $\theta$   
3:  $\alpha \leftarrow 1.25 \times \Delta/n$   
4: Initialize  $\theta$ , let  $\theta' = \theta$  and  $u = 0$   
5: for  $t = 0$  to  $T - 1$  do  
6:   Sample  $(x, y)$  from  $\mathcal{Z}$   
7:   for  $i = 0$  to  $n - 1$  do  
8:      $x^{\text{adv}} \leftarrow \text{Proj}_{\{x': \|x' - x\|_\infty \leq \Delta\}} \{x +$   
         $\alpha \nabla_x f(\theta; x, y)\}$   
9:   end for  
10:   $\theta \leftarrow \text{Proj}_{\{\theta'': \|\theta''\|_0 \leq k\}} (\theta -$   
       $\gamma_t \nabla_\theta [f(\theta; x^{\text{adv}}, y) + \frac{\rho}{2} \|\theta - \theta' + u\|_F^2])$   
11:   $\theta' \leftarrow \text{ZeroKmeans}_{2^b}(\theta + u)$   
12:   $u \leftarrow u + (\theta - \theta')$   
13: end for
```

Empirical Study

- Datasets and Benchmark Models
 - Four popular image classification datasets
 - Pick one top-performer CNN model on each

Table 1: The datasets and CNN models used in the experiments.

Models	#Parameters	bit width	Model Size (bits)	Dataset & Accuracy
LeNet	430K	32	13,776,000	MNIST: 99.32%
ResNet34	21M	32	680,482,816	CIFAR-10: 93.67%
ResNet34	21M	32	681,957,376	CIFAR-100: 73.16%
WideResNet	11M	32	350,533,120	SVHN: 95.25%

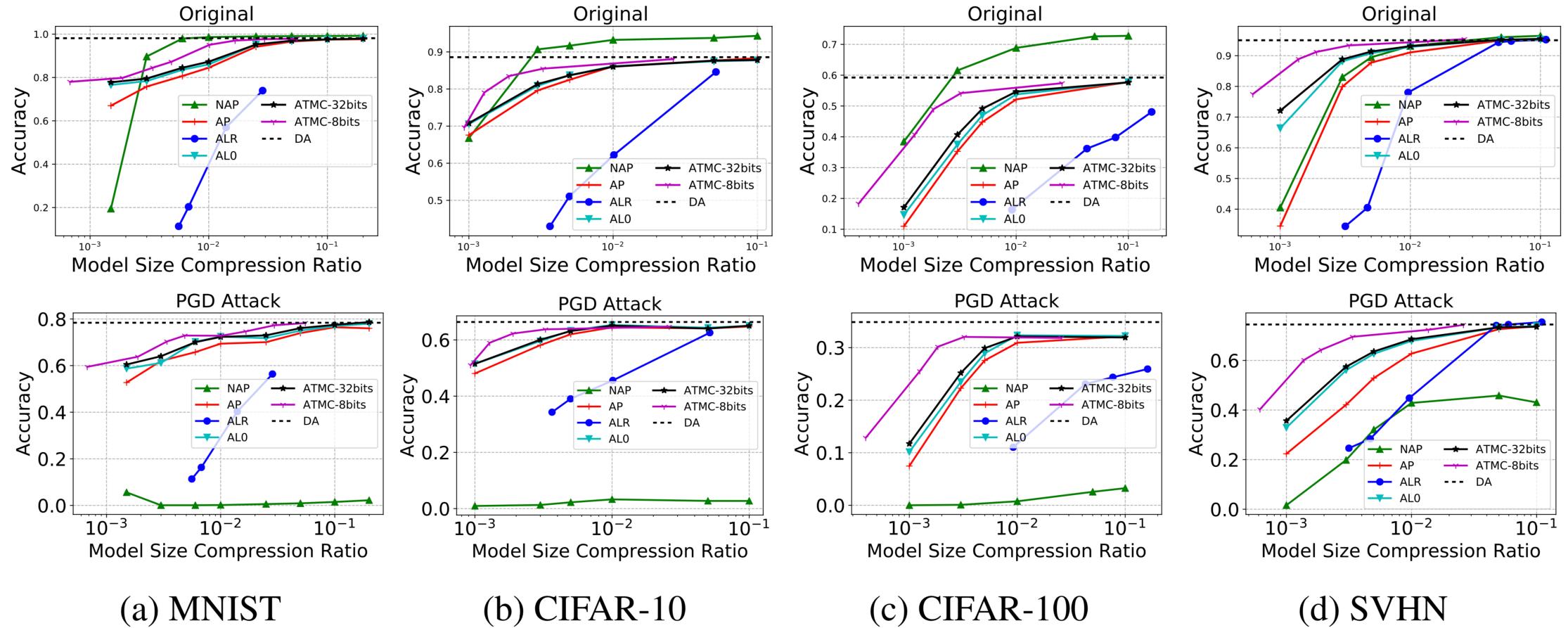
Empirical Study

- Evaluation Metrics
 - Classification accuracies on both benign and on attacked testing sets are reported
- Model Size
 - multiplying the quantization bit per element with the total number of non-zero elements, added with the storage size for the quantization thresholds
- Compression Ratio
 - the ratio between the compressed and original model sizes
- Adversarial Attack Settings
 - PGD attack with the same settings as used in adversarial training on testing sets to evaluate model robustness

Empirical Study

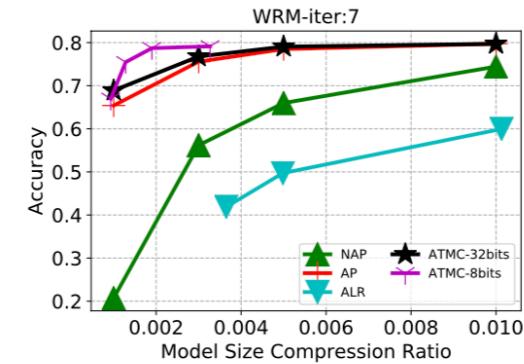
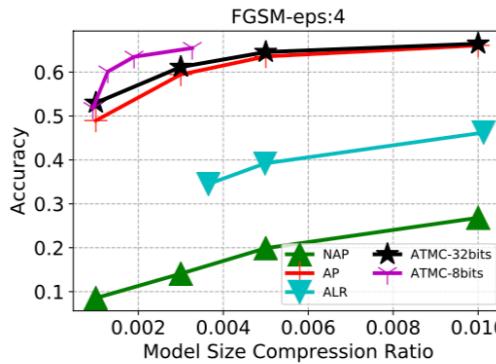
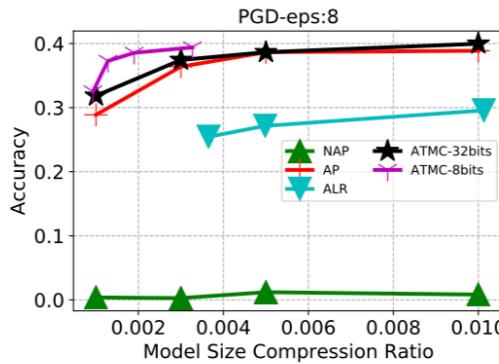
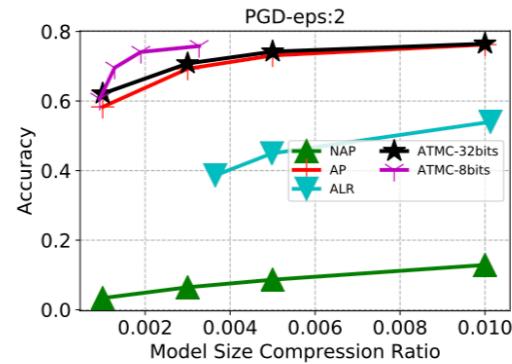
- Baselines
 - Non-adversarial Pruning (NAP)
 - Dense Adversarial Training (DA)
 - Adversarial Pruning (AP)
 - Adversarial I0 Pruning (AI0)
 - Adversarial Low-Rank Decomposition (ALR)
 - ATMC (8bits, 32 bits)
- Comparison on same compression ratio

Empirical Study



Empirical Study

- ATMC also maintains robustness under different adversarial attacks

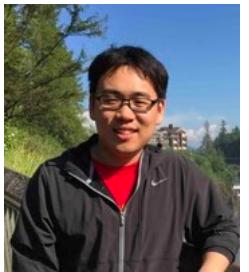


(a) PGD, perturbation=2 (b) PGD, perturbation=8 (c) FGSM, perturbation=4 (d) WRM, penalty=1.3, iteration=7

Conclusion

- We propose the ATMC framework, by integrating the two goals, model compression and robustness in one unified constrained optimization framework
- Our extensive experiments endorse the effectiveness of ATMC
 - Naïve model compression may hurt robustness, if the latter is not explicitly taken into account
 - A proper joint optimization could achieve both well
 - A properly compressed model could even maintain almost the same accuracy and robustness compared to the original one

Thanks



Shupeng Gui



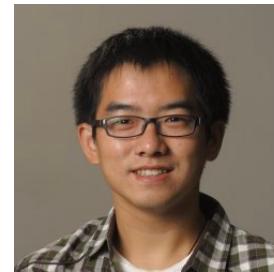
Haotao Wang



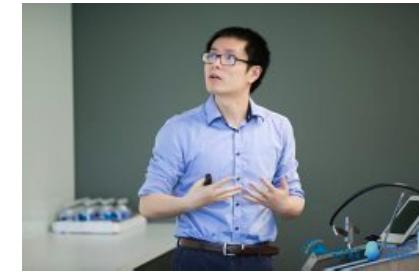
Haichuan Yang



Chen Yu



Zhangyang Wang



Ji Liu



VITA



Reference

- Guo, Yiwen, et al. "Sparse dnns with improved adversarial robustness." *Advances in neural information processing systems*. 2018.
- Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." *arXiv preprint arXiv:1412.6572* (2014).
- Madry, Aleksander, et al. "Towards deep learning models resistant to adversarial attacks." *arXiv preprint arXiv:1706.06083* (2017).
- Han, Song, et al. "Learning both weights and connections for efficient neural network." *Advances in neural information processing systems*. 2015.
- Xu, Huan, Constantine Caramanis, and Shie Mannor. "Sparse algorithms are not stable: A no-free-lunch theorem." *IEEE transactions on pattern analysis and machine intelligence* 34.1 (2011): 187-193.
- Tsipras, Dimitris, et al. "Robustness may be at odds with accuracy." *arXiv preprint arXiv:1805.12152* (2018).
- Li, Hao, et al. "Pruning filters for efficient convnets." *arXiv preprint arXiv:1608.08710* (2016).
- Lloyd, "Least squares quantization in pcm", IEEE transactions on information theory, 1982
- Han, Song, Huizi Mao, and William J. Dally. "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding." *arXiv preprint arXiv:1510.00149* (2015).
- He, Yihui, Xiangyu Zhang, and Jian Sun. "Channel pruning for accelerating very deep neural networks." *Proceedings of the IEEE International Conference on Computer Vision*. 2017.
- Zhou, Hao, Jose M. Alvarez, and Fatih Porikli. "Less is more: Towards compact cnns." *European Conference on Computer Vision*. Springer, Cham, 2016.