

# **Fine-Grained Iterative Adversarial Attacks with Limited Computation Budget**

Zhichao Hou, Weizhi Gao, Xiaorui Liu

# Attacks Are Essential Yet Expensive

## Attack perspective:

- Iterative attack such as PGD (Madry et al., 2017) are widely adopted as strong oracles to benchmark the robustness.
- PGD with 20, 50, or even 100 steps is typically used, which incurs roughly **40–200×** the cost of a natural inference.

## Defense perspective:

- The most effective defense: adversarial training (Madry et al., 2017; Zhang et al., 2019; Wang et al., 2019) critically depends on strong iterative attacks to generate challenging adversaries during training.
- The standard PGD-10 procedure already costs around **10×** more than natural training.

**Critical Challenge:** *Given a prescribed computation budget, how can we maximize the strength of iterative adversarial attacks ?*

# Expensive Computation in Attacks

---

**Algorithm 1** Gradient-based Adversarial Attack

---

**Require:** clean input  $\mathbf{x}$ , label  $y$ , adversarial loss  $\mathcal{L}$ , attack budget  $\epsilon$ , step size  $\alpha$ , iteration  $T$ .

**Ensure:** adversarial example  $\tilde{\mathbf{x}}$

1:  $\mathbf{x}_1 \leftarrow \mathbf{x}$  ▷ initialization

2: **for**  $t = 1, \dots, T$  **do**

3:  $\mathbf{x}_t \rightarrow \mathbf{o}_t^{(1)} \dots \mathbf{o}_t^{(l)} \dots \rightarrow \mathbf{o}_t^{(L)} \rightarrow \mathcal{L}$

4:  $\frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} \leftarrow \frac{\partial \mathcal{L}}{\partial \mathbf{o}_t^{(1)}} \dots \frac{\partial \mathcal{L}}{\partial \mathbf{o}_t^{(l)}} \dots \leftarrow \frac{\partial \mathcal{L}}{\partial \mathbf{o}_t^{(L)}} \leftarrow \mathcal{L}$

5:  $\mathbf{x}_{t+1} \leftarrow \Pi_{\mathcal{B}_\epsilon(\mathbf{x})} \left( \mathbf{x}_t + \alpha \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} \right)$

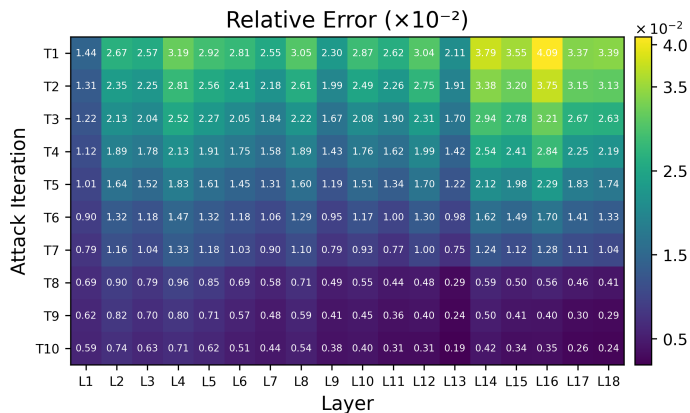
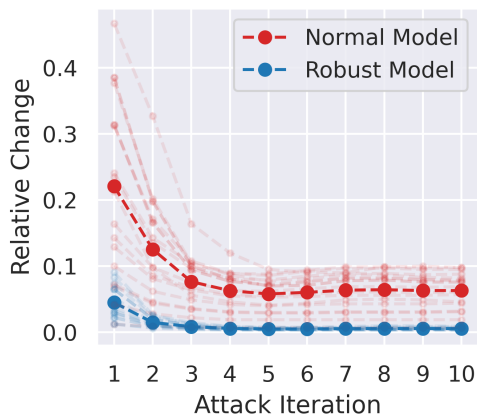
6: **end for**

7: **return**  $\tilde{\mathbf{x}} \leftarrow \mathbf{x}_{T+1}$

---

Gradient-based iterative attacks are notoriously expensive because it requires **multiple iterations**, and each step requires **both a forward and backward pass**.

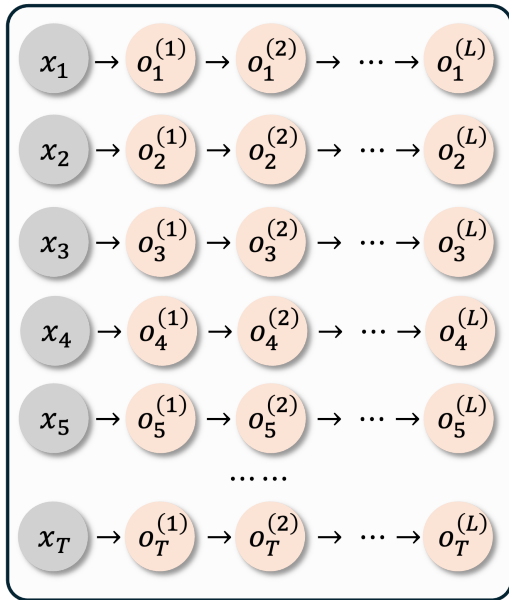
# Computation Redundancy in Attacks



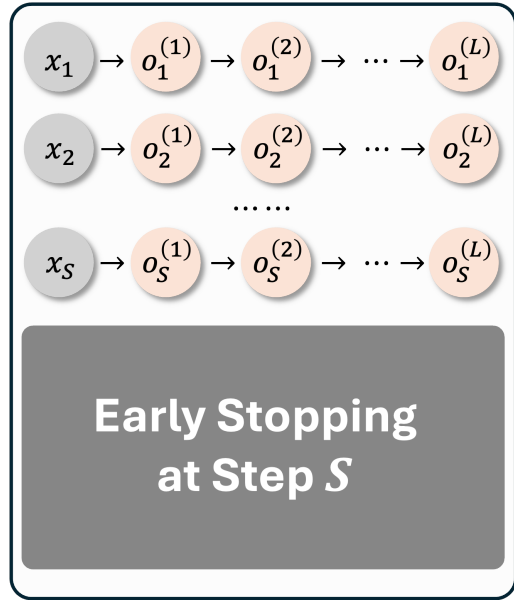
Activation relative change  $\|\mathbf{a}_t - \mathbf{a}_{t-1}\|/\|\mathbf{a}_t\|$  for ResNet-18 on CIFAR-10.

- Figure (Left) shows that the activations become highly similar after a small number of iterations. This indicates substantial redundancy in repeated computations across attack iterations. Additionally, the “Robust Model” exhibits markedly smaller activation changes than the “Normal Model”;
- Figure (Right) shows that while all layers follow the same overall decreasing trend, the decay rate of relative activation change varies across layer.

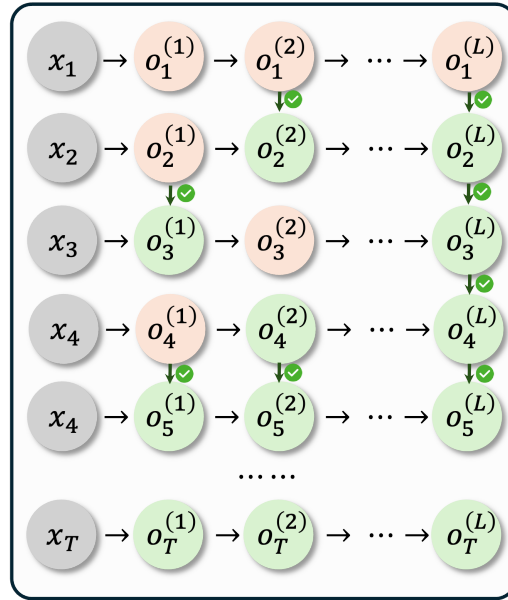
# Spiking Iterative Attack



(a) Vanilla Iterative Attack



(b) Coarse-Grained Attack



(c) Fine-Grained Attack

# Attack as Combinatorial Optimization

Let  $\Delta = (\delta_{t,l}) \in \{0, 1\}^{T \times L}$  be a binary mask where  $\delta_{t,l} = 1$  indicates that layer  $l$  is fully computed at attack iteration  $t$ , while  $\delta_{t,l} = 0$  indicates reuse of previously computed activations. The corresponding *fine-grained* optimization is the following:

$$\max_{\Delta \in \{0,1\}^{T \times L}} \mathcal{L}(\mathbf{x}_{T+1}(\Delta), y) \quad \text{s.t.} \quad \sum_{t=1}^T \sum_{l=1}^L C_{t,l} \delta_{t,l} \leq C_{\text{total}}, \quad (2)$$

where  $\mathbf{x}_{T+1}(\Delta)$  denotes the final perturbed input after  $T$  attack steps executed under mask  $\Delta$ .

## Challenges:

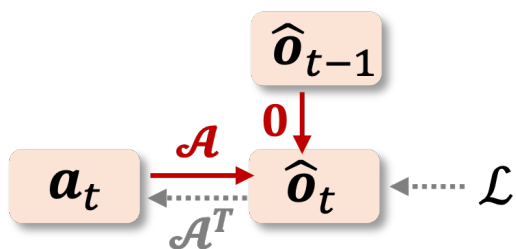
- (1) Vast search space:  $\mathcal{O}(2^{T \times L})$ ;
- (2) Expensive objective evaluations;
- (3) Broken gradient flow;

# Spiking Forward Computation

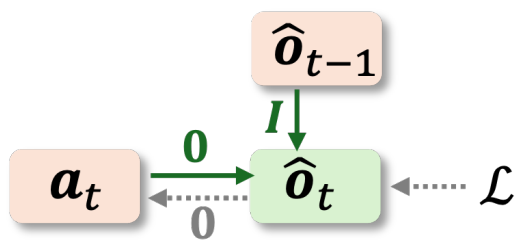
$$\left\{ \begin{array}{l} \mathbf{o}_1^{(l)} = \mathcal{A}^{(l)}(\mathbf{a}_1^{(l)}) \\ \mathbf{o}_2^{(l)} = \mathcal{A}^{(l)}(\mathbf{a}_2^{(l)}) = \mathcal{A}^{(l)}(\mathbf{a}_2^{(l)} - \mathbf{a}_1^{(l)}) + \mathbf{o}_1^{(l)} \\ \dots \\ \mathbf{o}_t^{(l)} = \mathcal{A}^{(l)}(\mathbf{a}_t^{(l)}) = \mathcal{A}^{(l)}(\mathbf{a}_t^{(l)} - \mathbf{a}_{t-1}^{(l)}) + \mathbf{o}_{t-1}^{(l)} \\ \dots \\ \mathbf{o}_T^{(l)} = \mathcal{A}^{(l)}(\mathbf{a}_T^{(l)}) = \mathcal{A}^{(l)}(\mathbf{a}_T^{(l)} - \mathbf{a}_{T-1}^{(l)}) + \mathbf{o}_{T-1}^{(l)} \end{array} \right. \Rightarrow \left\{ \begin{array}{l} \mathbf{o}_1^{(l)} = \hat{\mathbf{o}}_1^{(l)} = \mathcal{A}^{(l)}(\mathbf{a}_1^{(l)}) \\ \mathbf{o}_2^{(l)} \approx \hat{\mathbf{o}}_2^{(l)} = \mathcal{A}^{(l)}(\mathcal{S}_\rho(\mathbf{a}_2^{(l)}, \mathbf{a}_1^{(l)})) + \hat{\mathbf{o}}_1^{(l)} \\ \dots \\ \mathbf{o}_t^{(l)} \approx \hat{\mathbf{o}}_t^{(l)} = \mathcal{A}^{(l)}(\mathcal{S}_\rho(\mathbf{a}_t^{(l)}, \mathbf{a}_{t-1}^{(l)})) + \hat{\mathbf{o}}_{t-1}^{(l)} \\ \dots \\ \mathbf{o}_T^{(l)} \approx \hat{\mathbf{o}}_T^{(l)} = \mathcal{A}^{(l)}(\mathcal{S}_\rho(\mathbf{a}_T^{(l)}, \mathbf{a}_{T-1}^{(l)})) + \hat{\mathbf{o}}_{T-1}^{(l)} \end{array} \right.$$

$$\mathcal{S}_\rho(\mathbf{a}_t, \mathbf{a}_{t-1}) = \begin{cases} \mathbf{1} \cdot (\mathbf{a}_t - \mathbf{a}_{t-1}), & \|\mathbf{a}_t - \mathbf{a}_{t-1}\| / \|\mathbf{a}_t\| \geq \rho, \\ \mathbf{0} \cdot (\mathbf{a}_t - \mathbf{a}_{t-1}), & \|\mathbf{a}_t - \mathbf{a}_{t-1}\| / \|\mathbf{a}_t\| < \rho, \end{cases}$$

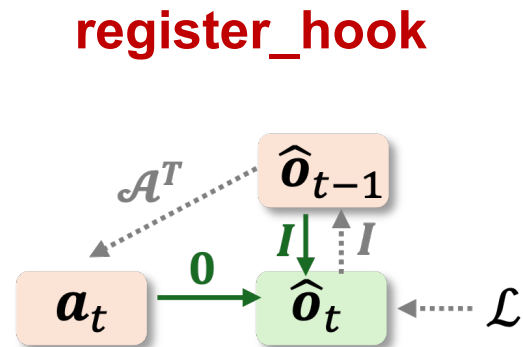
# Virtual Backward Computation



(a) Exact Gradient



(b) Vanishing Gradient



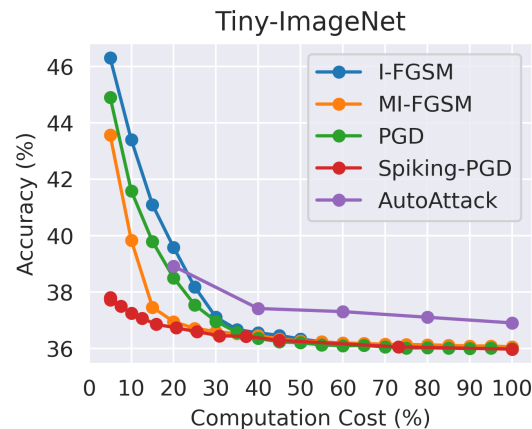
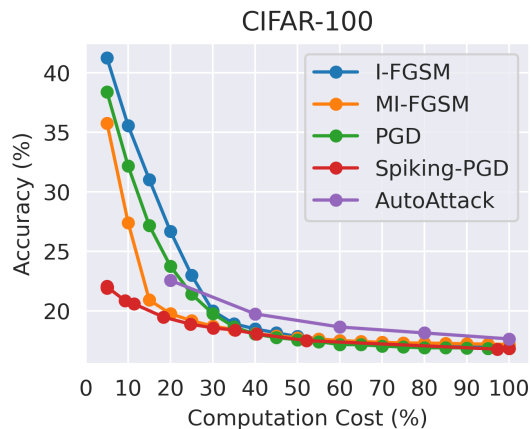
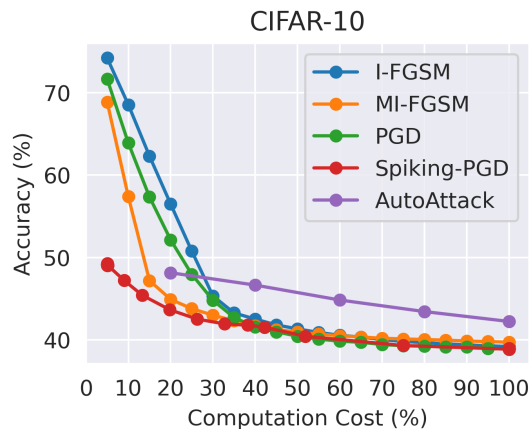
(c) Virtual Gradient

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}_t^{(l)}} = \mathcal{S}'_{\rho} \cdot \mathcal{A}^{(l)\top} \left( \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{o}}_t^{(l)}} \right) = \begin{cases} \mathcal{A}^{(l)\top} \left( \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{o}}_t^{(l)}} \right), & \|\mathbf{a}_t - \mathbf{a}_{t-1}\| / \|\mathbf{a}_t\| \geq \rho, \Rightarrow \text{Figure 4 (a)} \\ \mathbf{0}, & \|\mathbf{a}_t - \mathbf{a}_{t-1}\| / \|\mathbf{a}_t\| < \rho, \Rightarrow \text{Figure 4 (b)} \end{cases}$$

# Adversarial Attack Strength



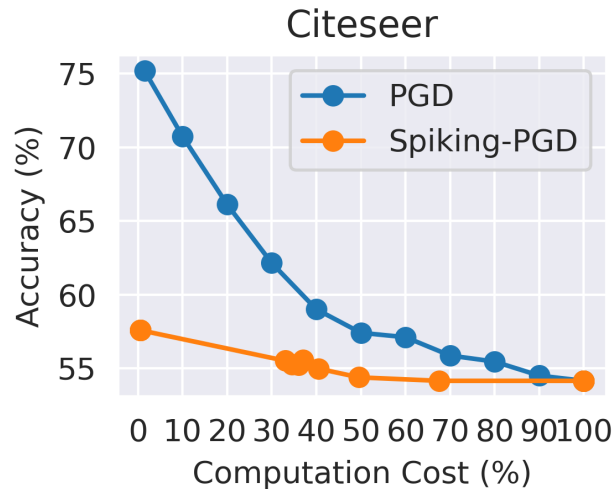
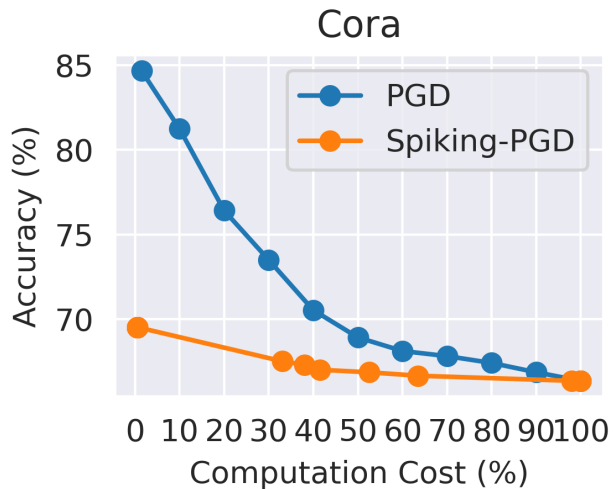
## Pixel Perturbation on Images



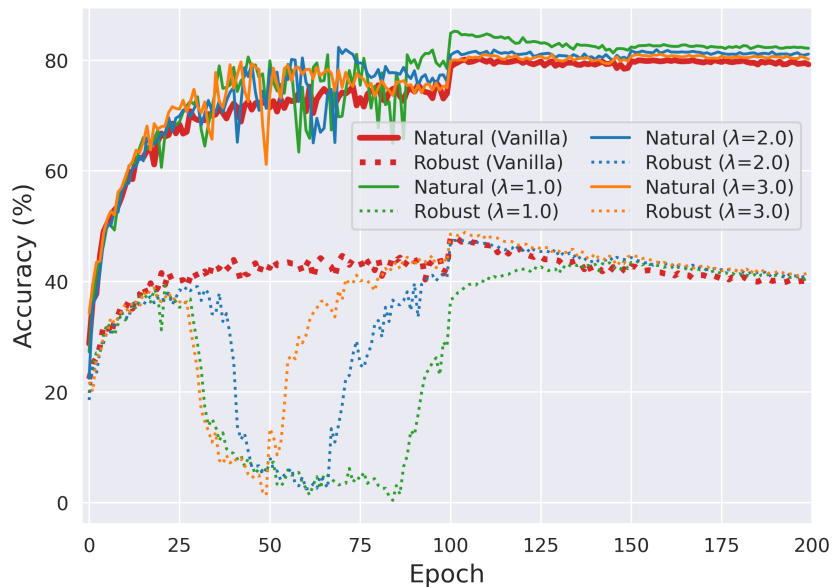
# Adversarial Attack Strength



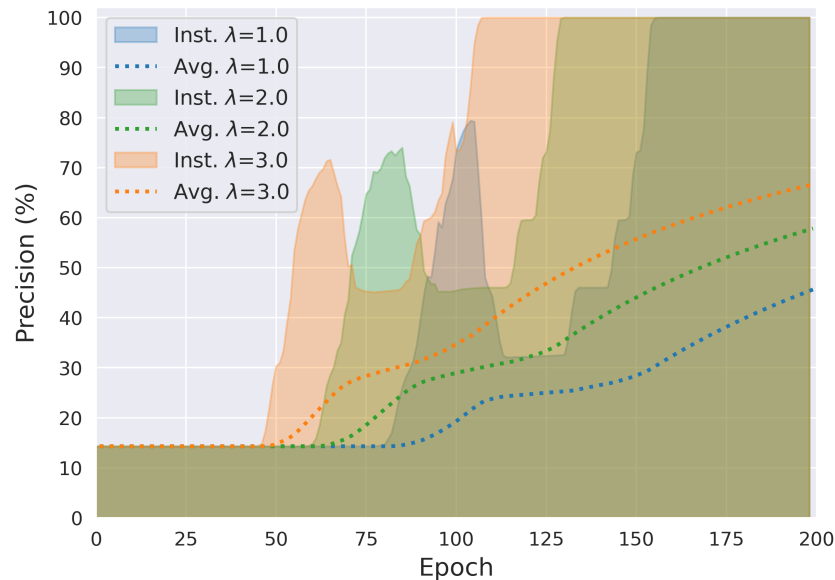
## Structure Attack on Graphs



# Efficient Adversarial Training

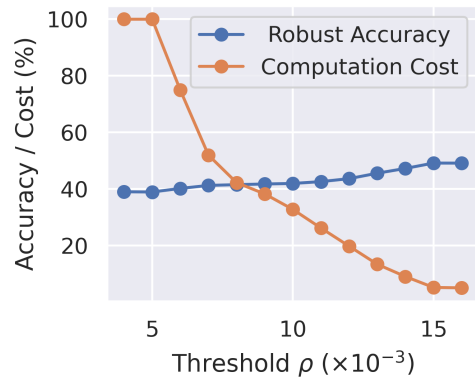


Accuracy curves of adversarial training on CIFAR10.

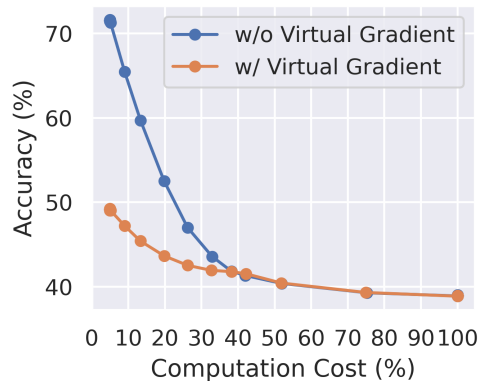


Precision curves of adversarial training on CIFAR-10.

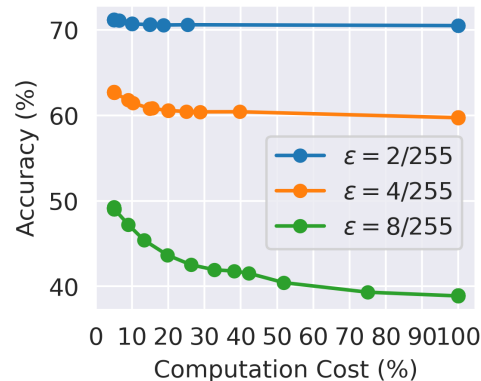
# Ablation Study



(a) Ablation on threshold  $\rho$



(b) Ablation on virtual gradient.



(c) Ablation on attack radius.



**NC STATE**  
UNIVERSITY

# Thanks for Listening

Zhichao Hou, Weizhi Gao, Xiaorui Liu