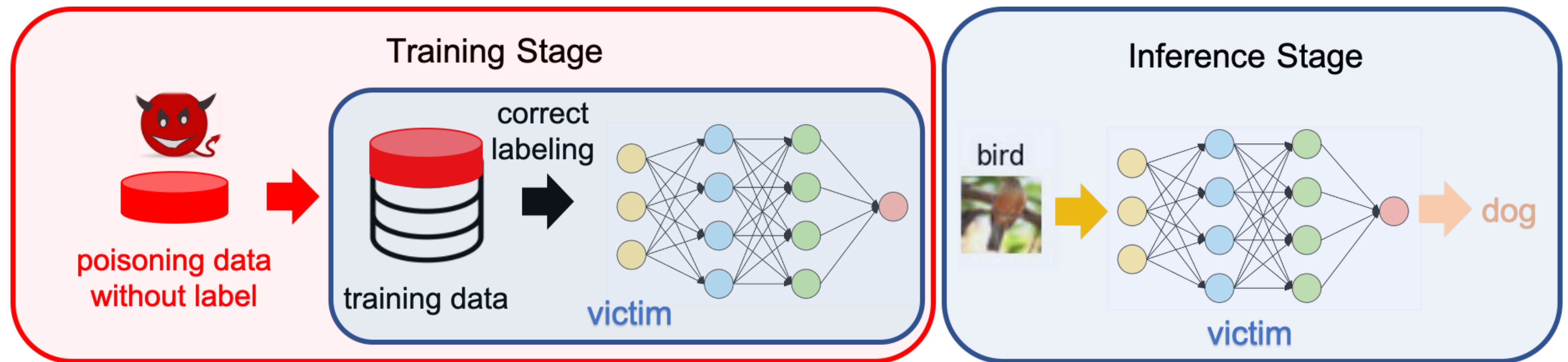


# First-Order Efficient General-Purpose Clean-Label Data Poisoning

Tianhang Zheng, Baochun Li  
University of Toronto

# Clean-label Data Poisoning

- Compromise Deep Learning at the Training Stage
- Insert poisoning data ( $< 10\%$ ) into training set without control on the labeling process



cannot insert the target data with wrong labels (will be corrected)

# Related Work: Influence Function

- Influence Function: characterize the effect that  $\mathbf{x} \rightarrow \mathbf{x} + \delta$  has on the loss of a target (testing) sample  $\mathbf{x}_{test}$

$$[-\nabla_{\Theta} \mathcal{L}(\mathbf{x}_{test}, y_{test}, \Theta)^T \underbrace{\mathbf{H}_{\Theta}^{-1} \nabla_{\mathbf{x}} \nabla_{\Theta} \mathcal{L}(\mathbf{x}, y, \Theta)}_{\text{Second-order terms}}] \cdot \delta \quad (1)$$

Second-order terms  $\mathbf{H}_{\Theta} = \frac{1}{N} \sum_{n=1}^N \nabla_{\Theta}^2 \mathcal{L}(\mathbf{x}_n, y_n, \Theta)$

- To decrease the loss, we can set

$$\delta = \epsilon [\nabla_{\Theta} \mathcal{L}(\mathbf{x}_{test}, y_{test}, \Theta)^T \mathbf{H}_{\Theta}^{-1} \nabla_{\mathbf{x}} \nabla_{\Theta} \mathcal{L}(\mathbf{x}, y, \Theta)]^T \quad (2)$$

# Related Work: Meta Poison

- The optimization process can be formulated as

$$\Theta_1 = \Theta - \gamma \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_p \cup \mathbf{X}_c), Y_p \cup Y_c)}{\partial \Theta}$$

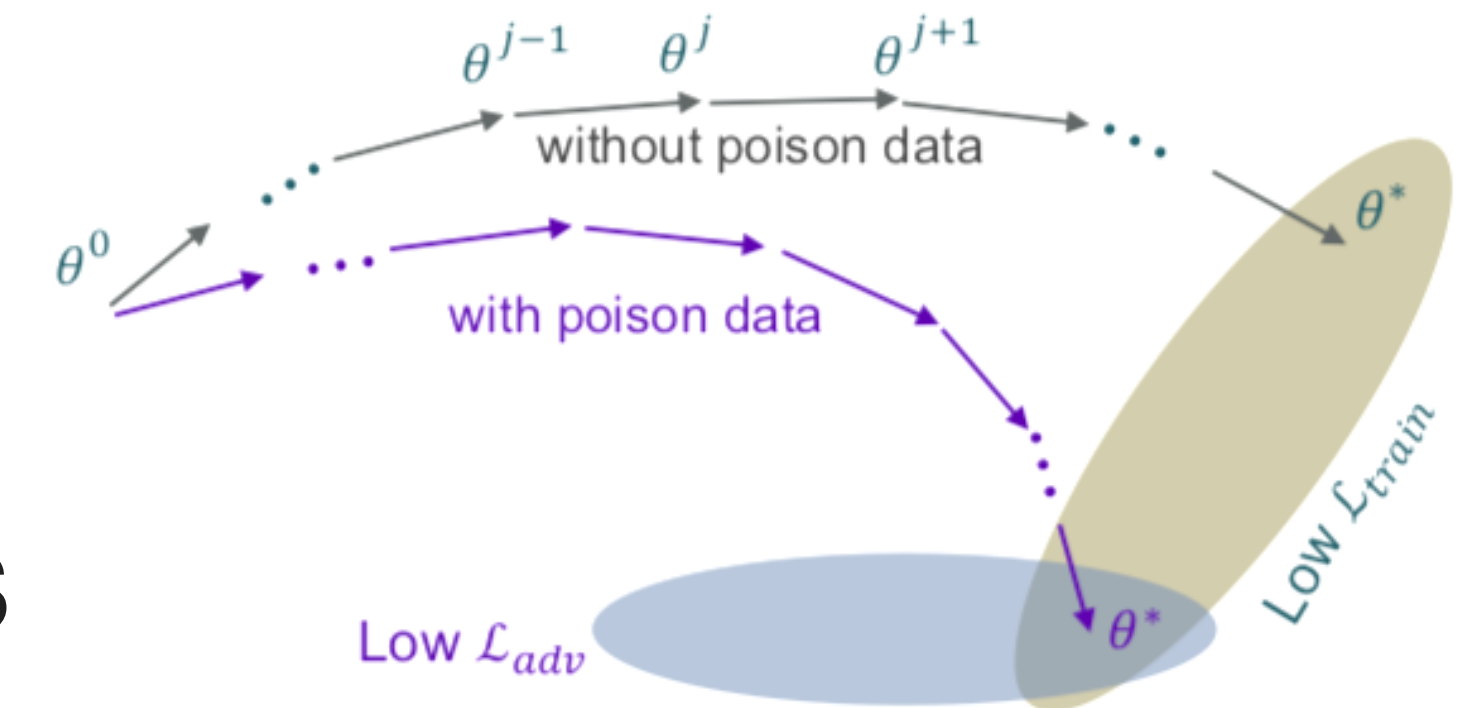
Model weight updates on the current dataset

$$\Theta_2 = \Theta_1 - \gamma \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta_1}(\mathbf{X}_p \cup \mathbf{X}_c), Y_p \cup Y_c)}{\partial \Theta_1}$$

$$\mathbf{X}_p = \mathbf{X}_p - \beta \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta_2}(\mathbf{X}_t), Y_t)}{\partial \mathbf{X}_p} \text{ (update } \mathbf{X}_p \text{)}$$

Update the poisoning subset

$$\frac{\partial^2 \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_p \cup \mathbf{X}_c), Y_p \cup Y_c)}{\partial \Theta \partial \mathbf{X}_p} \text{ to compute } \frac{\partial \Theta_1}{\partial \mathbf{X}_p}$$

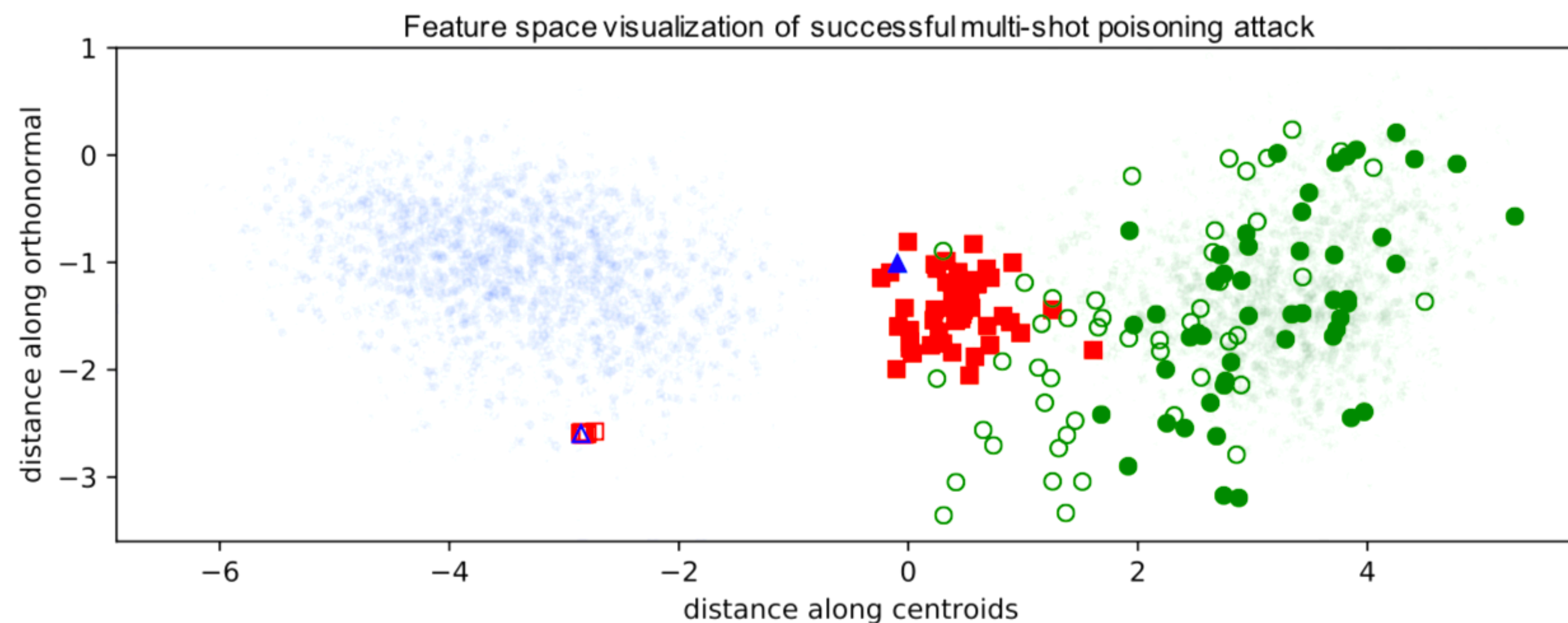


# Related Work: Feature Collision

- Search for feature collision between  $\mathbf{x} + \boldsymbol{\delta}$  (from the attack targeted class) and  $\mathbf{x}_{target}$

$$\operatorname{argmin}_{\boldsymbol{\delta}} \|\mathbf{f}_{\boldsymbol{\Theta}}(\mathbf{x} + \boldsymbol{\delta}) - \mathbf{f}_{\boldsymbol{\Theta}}(\mathbf{x}_{target})\|_2^2 + \beta \|\boldsymbol{\delta}\|_2^2$$

mainly effective in transfer learning scenario  $\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{f}_{\boldsymbol{\Theta}}(\cdot))$





# First-order Poisoning Attack

- First-order refers to only using first-order derivative information
- Methodology (summary):
  - Identify the **first-order adversary-desired model update** that can push the model towards predicting the target data as the attack targeted label
  - Formulate **a necessary condition** to optimize the poisoning data
  - We prove that our first-order poisoning method is an approximation of a second-order approach with **theoretically-guaranteed performance**

# Desired Model Update

- An **adversary-desired model update** pushes the model towards recognizing the target data sample as attack targeted label

$$\delta_{\Theta} = \tilde{\Theta} - \Theta = -\alpha \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_t), Y_t)}{\partial \Theta} \quad (1)$$

- With the above update, the loss will be perturbed by a non-positive term

$$\frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_t), Y_t)}{\partial \Theta} \cdot \delta_{\Theta} = -\alpha \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_t), Y_t)}{\partial \Theta} \cdot \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_t), Y_t)}{\partial \Theta} \quad (2)$$

# Necessary Condition

- A necessary condition is formulated based on the desired model update and other first-order information

$$\mathcal{L}(\mathbf{F}_{\Theta+\delta_{\Theta}}(\mathbf{X}_p), Y_p) + \delta_p \cdot \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta+\delta_{\Theta}}(\mathbf{X}_p), Y_p)}{\partial \mathbf{X}_p} \leq \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_p), Y_p) + \delta_p \cdot \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_p), Y_p)}{\partial \mathbf{X}_p} \quad (1)$$

- Since we want  $\Theta + \delta_{\Theta} \approx \underset{\tilde{\Theta}}{\operatorname{argmin}} \mathcal{L}(\mathbf{X}_p + \delta_p, Y_p, \tilde{\Theta})$
- Necessary condition on  $\delta_p$  (with same direction of the opposite of the red part)

$$\delta_p \cdot \left( \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta+\delta_{\Theta}}(\mathbf{X}_p), Y_p)}{\partial \mathbf{X}_p} - \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_p), Y_p)}{\partial \mathbf{X}_p} \right) \ll 0 \quad (2)$$



# Theoretically Guaranteed Performance

- Our derived first-order update is an approximation of a second order poisoning update
  - The additional loss change caused by  $\delta_p$  can be approximated by

$$-\gamma \frac{\partial^2 \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_p), Y_p)}{\partial \Theta \partial \mathbf{X}_p} \cdot \delta_p \cdot \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_t), Y_t)}{\partial \Theta} \quad (1)$$

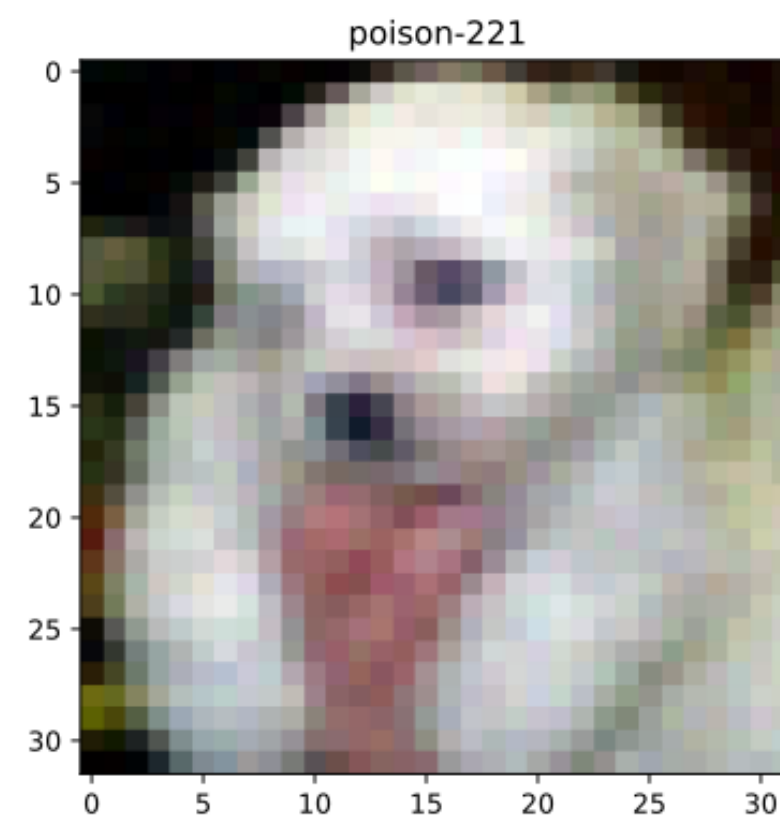
- The additional loss change is non-positive if

$$\delta_p = \epsilon \frac{\partial^2 \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_p), Y_p)}{\partial \mathbf{X}_p \partial \Theta} \cdot \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_t), Y_t)}{\partial \Theta} \quad (2)$$

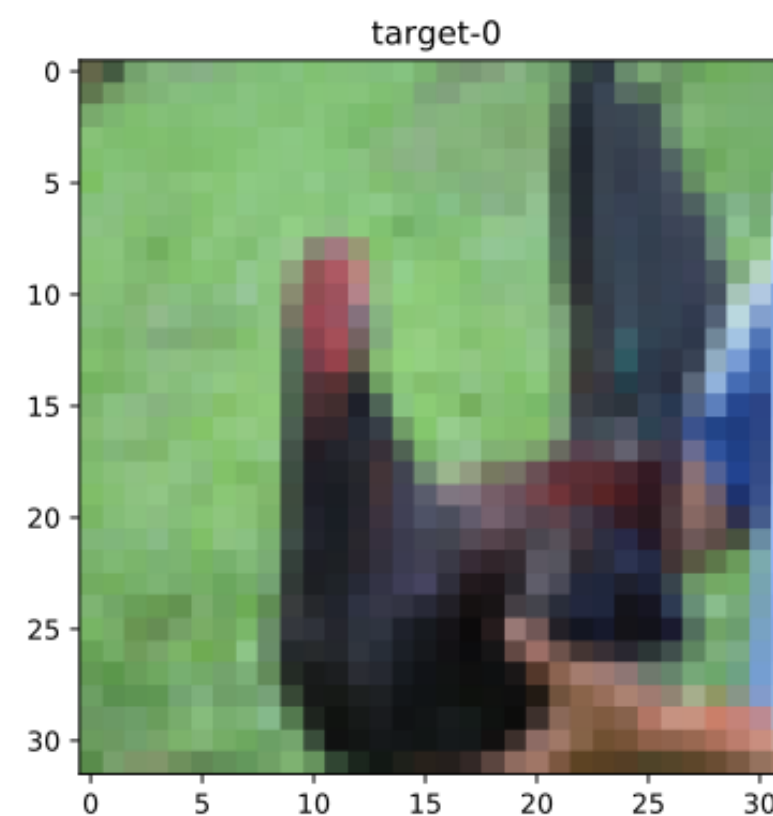
First-order update  $\delta_p = -\beta \left( \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta+\delta_{\Theta}}(\mathbf{X}_p), Y_p)}{\partial \mathbf{X}_p} - \frac{\partial \mathcal{L}(\mathbf{F}_{\Theta}(\mathbf{X}_p), Y_p)}{\partial \mathbf{X}_p} \right)$

# Watermark Trick

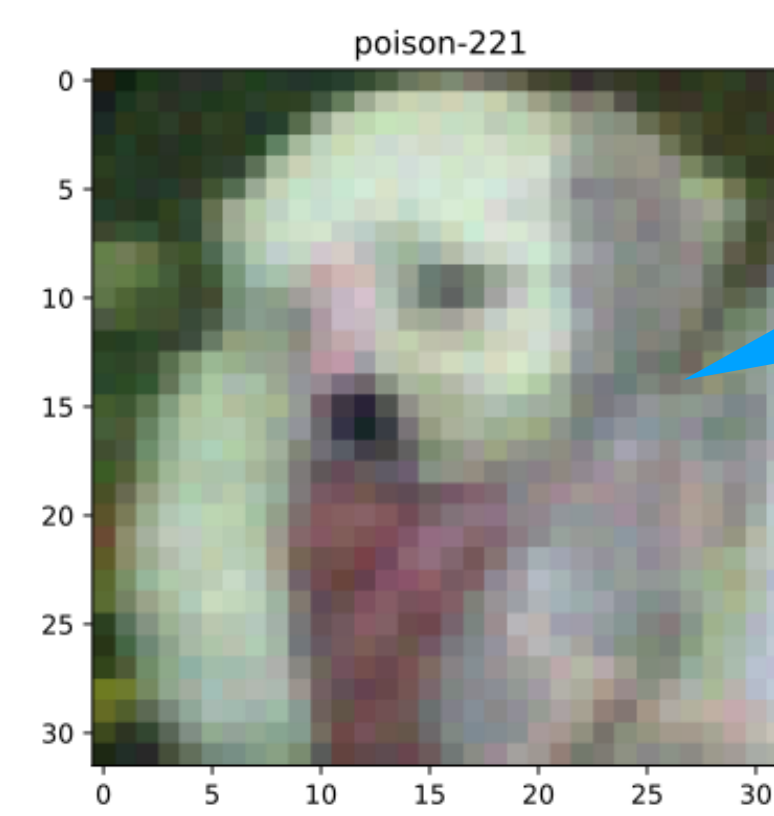
- ▶ Watermark:  $0.7 \times \text{Original Image} + 0.3 \times \text{Attack Target}$



Original image



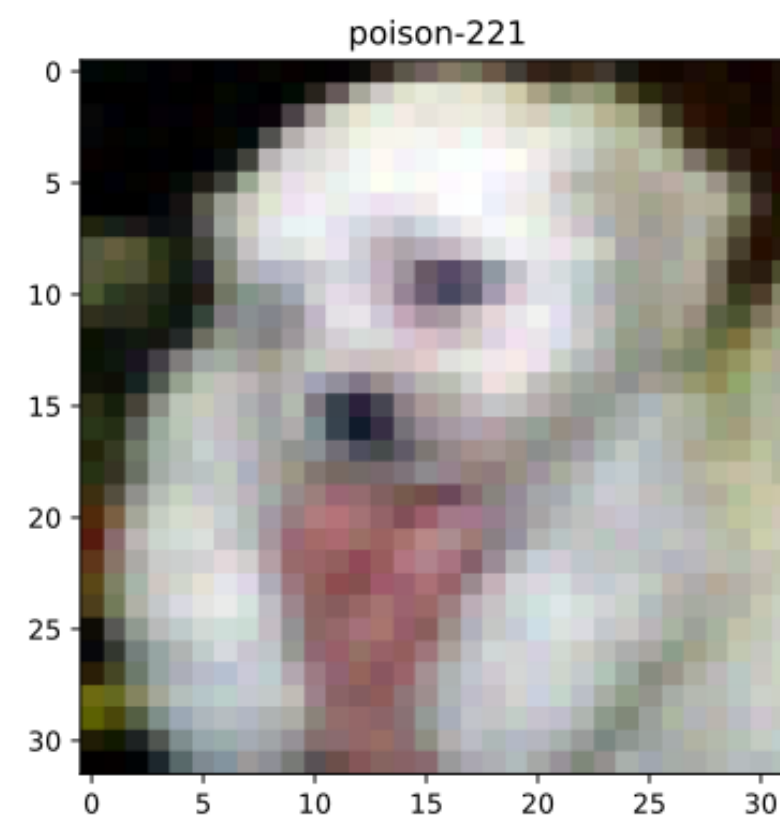
Attack target



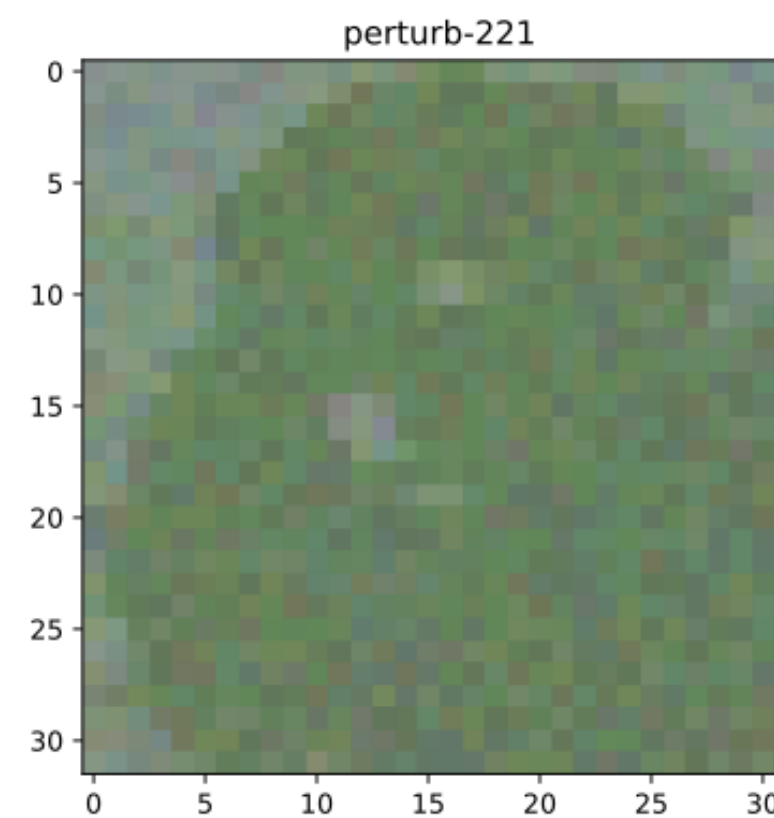
Watermarked image

# Color Perturbation

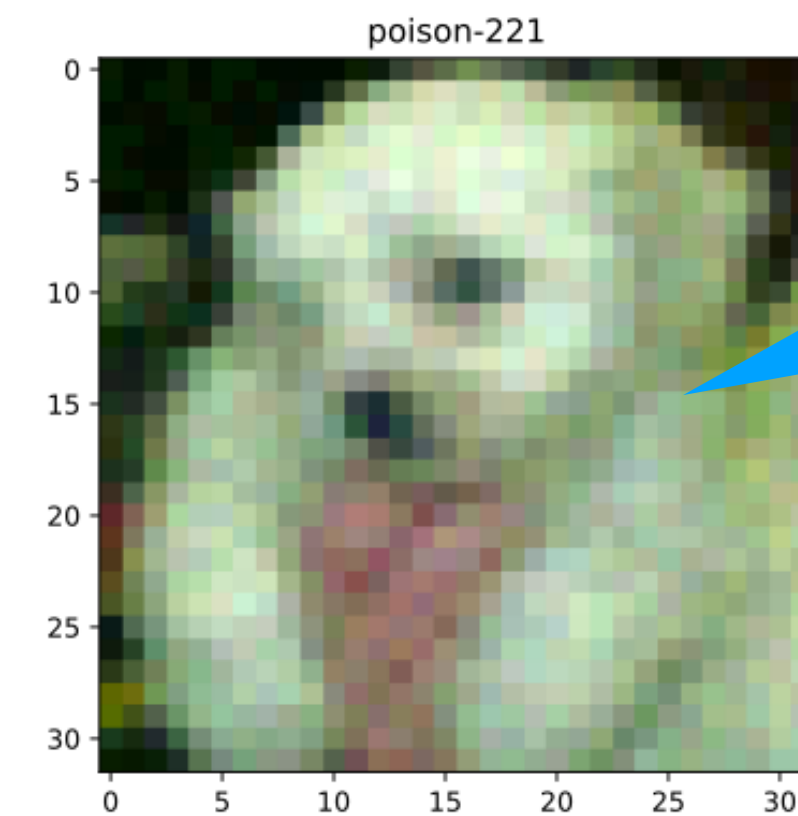
- Add color perturbation ( $\epsilon_c = 0.04$ ), crafted by meta-poison



Original image



Attack target

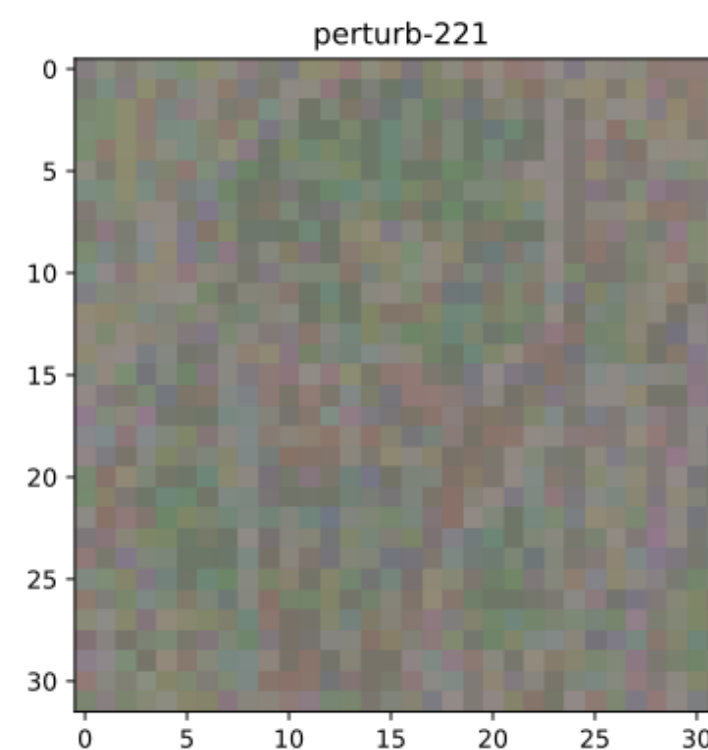
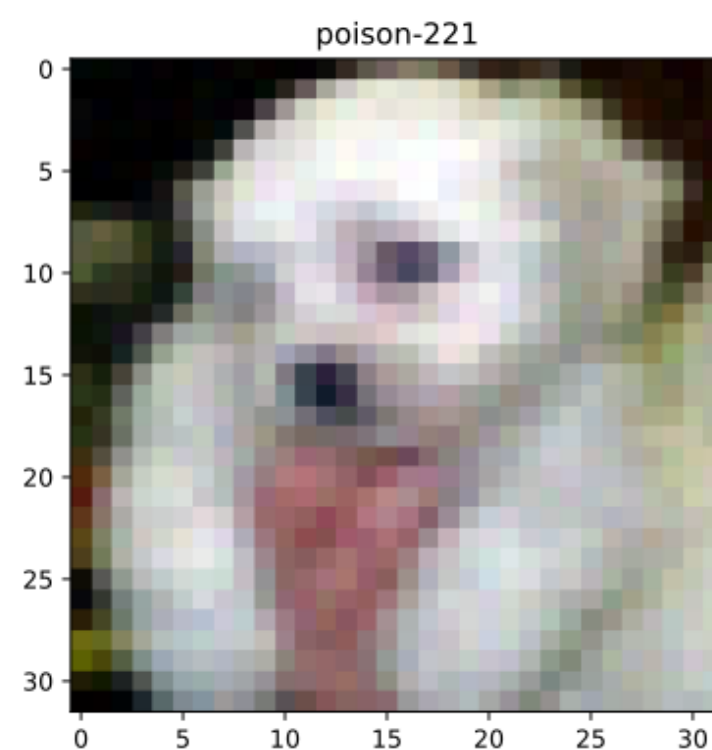


Watermarked image

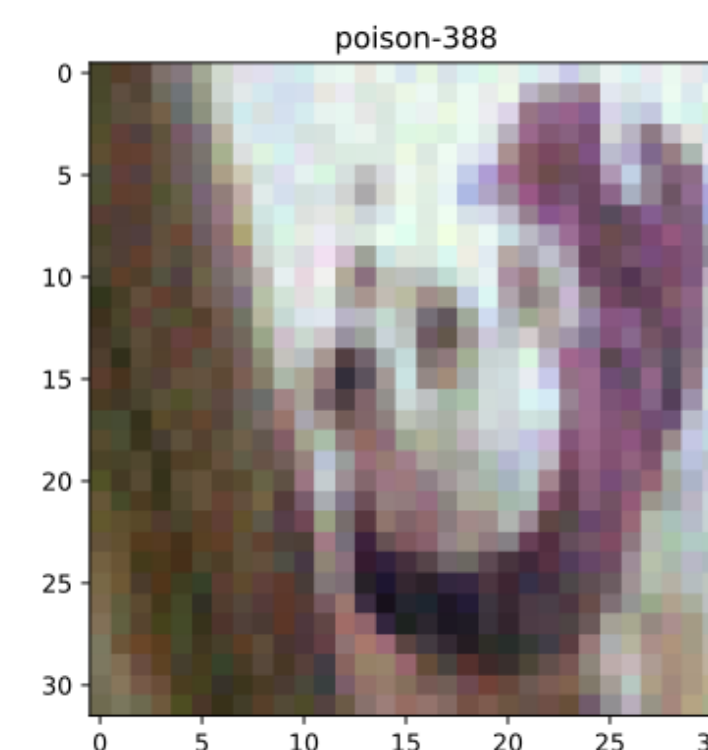
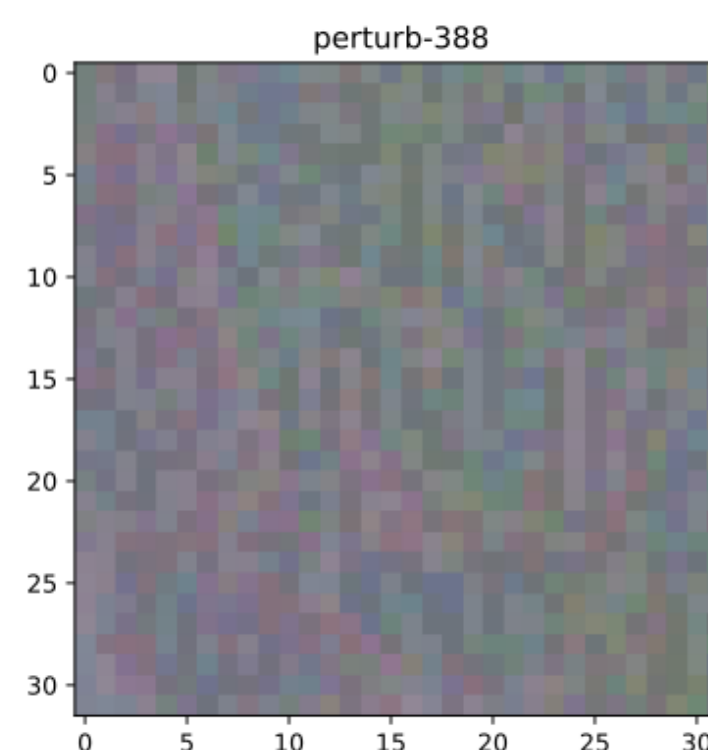
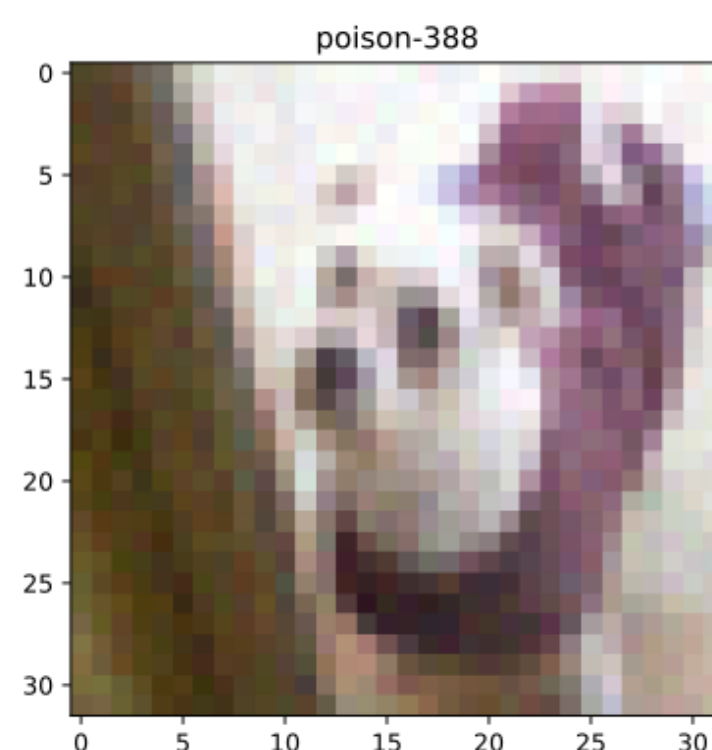
Green dog?

# Visualization

- Our setup: 10% watermark and  $\epsilon_c = 0.02$  color perturbation

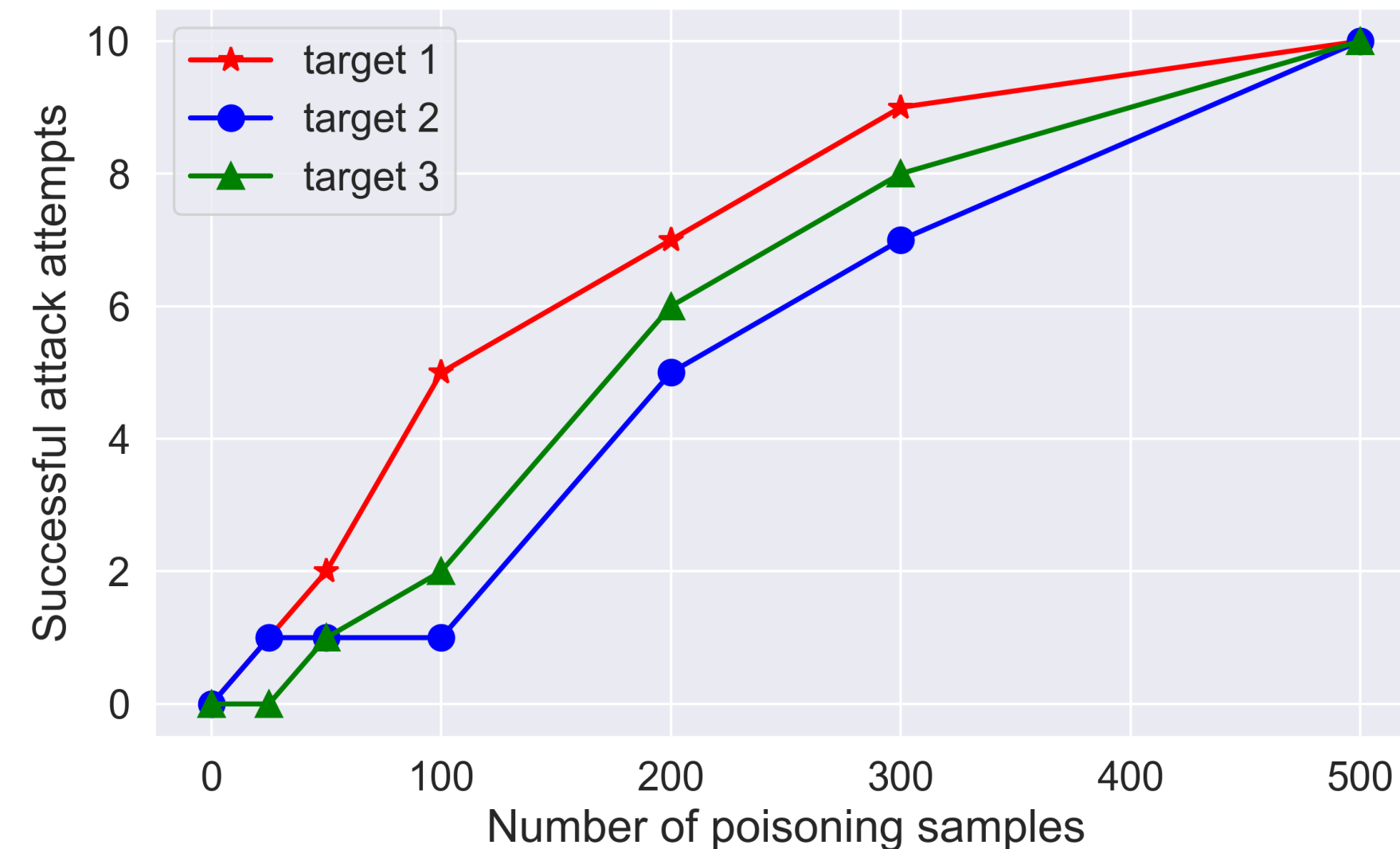


Look more  
normal



# Performance

- Our method is more efficient than meta-poison (NeurIPS20) and more general than feature collision (NeurIPS18)



The victim model is retrained from scratch, and our method can achieve 100% success rate with 500 poisoning samples here

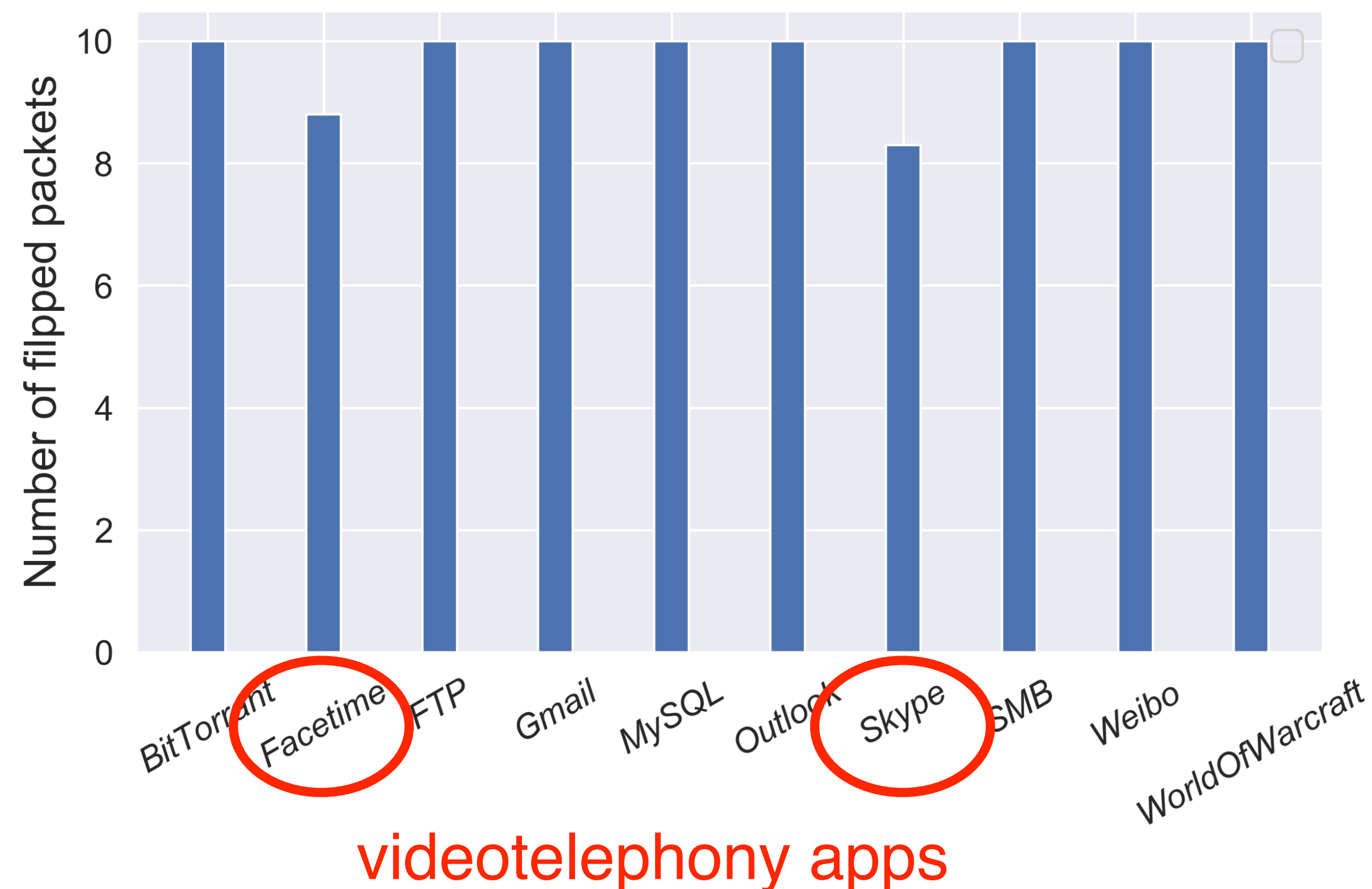
Model	# poison data	Meta Poison [22]	Our Method
ConvNetBN	50/50000	13.5s	9.3s
	500/50000	25.4s	16.0s
VGG-13	50/50000	35.5s	22.3s
	500/50000	71.7s	36.4s
ResNet-20	50/50000	48.6s	30.1s
	500/50000	95.4s	48.1s

Averaged computational time per crafting step on a single model on a single TITAN GPU



# Case Study: Network Traffic Classification

- The raw traffic data is processed into 2D images, and a CNN-based model is trained for classification



In each experiment, we select **10 target packets for each class**. We execute **10 attack attempts**, i.e., retrain 10 models on the poisoning data plus the remaining clean training data for each class.

# Conclusion Remarks

- First-Order Information Driven Clean-Label Data Poisoning:
  - More efficient than influence function and meta-poison
  - More general than feature collision
  - Theoretically guaranteed performance
- Future Work:
  - More efficient data poisoning
  - Defense against data poisoning

Contact: [th.zheng@mail.utoronto.ca](mailto:th.zheng@mail.utoronto.ca)