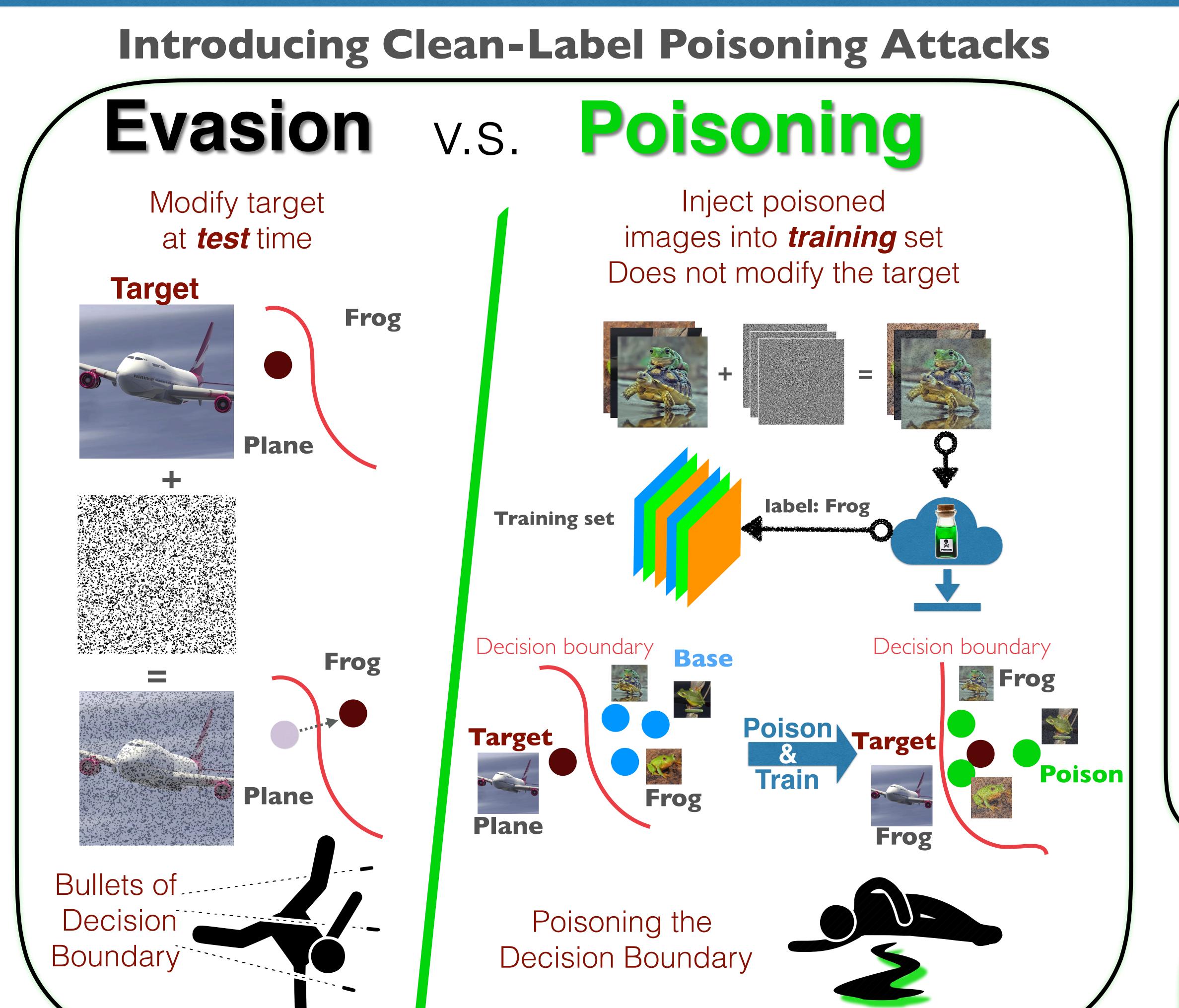


Transferable Clean-Label Poisoning Attacks on Deep Neural Nets

Chen Zhu*,1, W. Ronny Huang*,1, Ali Shafahi1, Hengduo Li1, Gavin Taylor2, Christoph Studer3, Tom Goldstein1 ¹University of Maryland, College Park, ²US Naval Academy, ³Cornell University



THREAT MODEL: CLEAN-LABEL ATTACKS Attacks can be executed by outsider WHY POISON? Poison data can be placed on the web You can't always control Poison data can be sent/emailed to data collectors target! Attacks are hard to detect for example... Performance only Clean-label: poisons are labeled changes on a Phishing/ selected target "correctly" security desk Competitor email HERA DE DE LES D

Crafting poisons

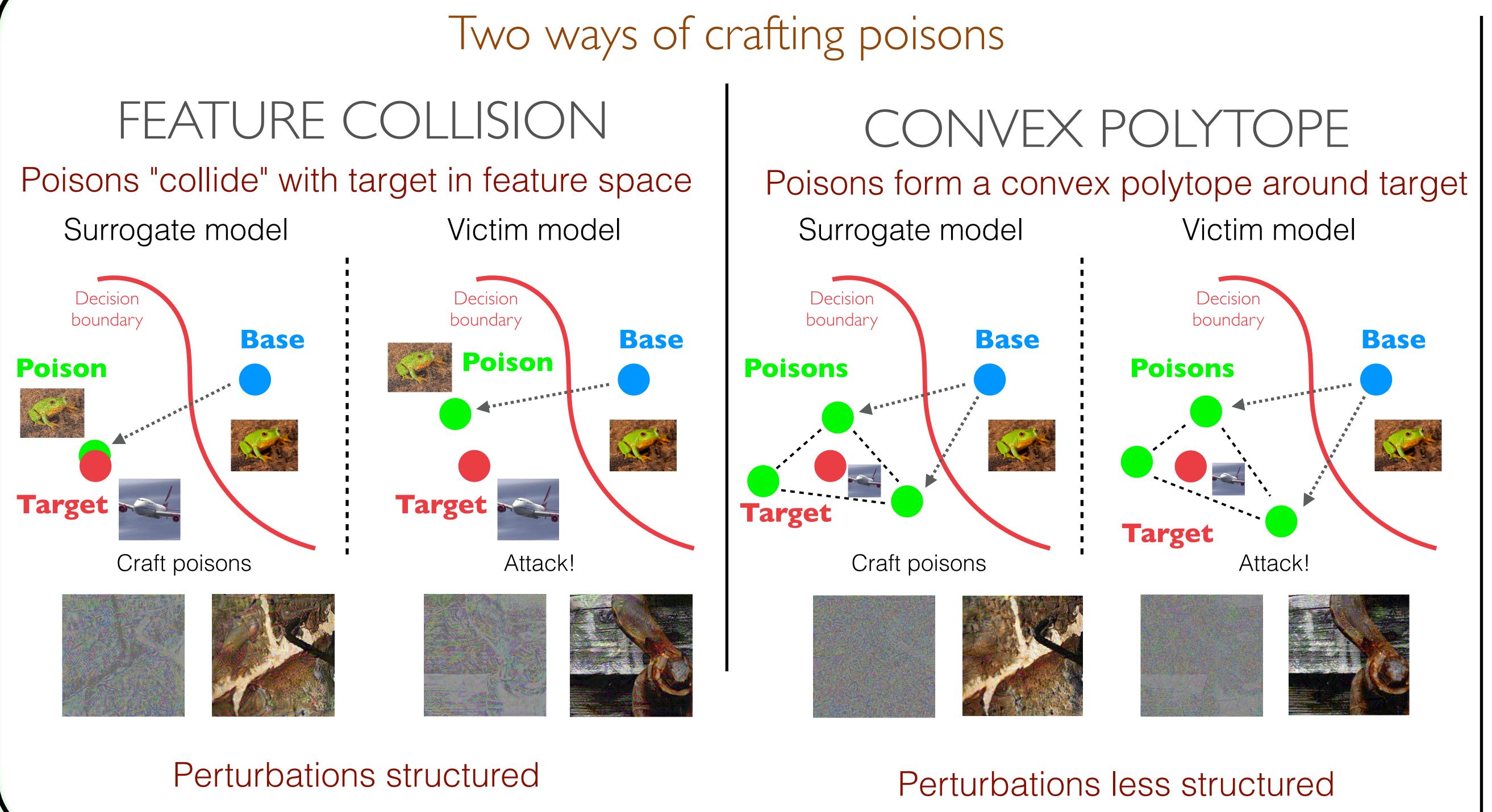
Implementation

bases / poisons

average over

m networks to make

poison transferable



OPTIMIZATION OBJECTIVE

Convex combination

of poisons

Close to target in

feature space

s.t. $\sum c_i^{(i)} = 1$ $c_i^{(i)} \ge 0$ $\forall i,j$ Enforce convex

Normalization

factor

combination

Dropout improves transferability Training Loss with Dropout Loss in Victim with Dropout

Algorithm 1 Convex Polytope Attack

Data: Clean base images $\{x_b^{(j)}\}_{j=1}^k$, substitute networks $\{\phi^{(i)}\}_{i=1}^m$, and maximum perturbation ϵ .

Result: A set of perturbed poison images $\{x_p^{(j)}\}_{j=1}^k$. Initialize $\boldsymbol{c}^{(i)} \leftarrow \frac{1}{k} \boldsymbol{1}, \boldsymbol{x}_p^{(j)} \leftarrow \boldsymbol{x}_b^{(j)}$

while not converged do for $i=1,\ldots,m$ do

 $A \leftarrow [\phi^{(i)}(\boldsymbol{x}_p^{(1)}), \dots, \phi^{(i)}(\boldsymbol{x}_p^{(k)})]$ $\alpha \leftarrow 1/\|A^{\mathsf{T}}A\|_2$ while not converged do $\boldsymbol{c}^{(i)} \leftarrow \boldsymbol{c}^{(i)} - \alpha A^{\top} (A \boldsymbol{c}^{(i)} - \phi^{(i)} (\boldsymbol{x}_t))$ project $c^{(i)}$ onto probability simplex

Gradient step on $\boldsymbol{x}_p^{(j)}$ with Adam

Clip $x_p^{(j)}$ so that the infinity norm constraint is satisfied.

Results



Pre-trained on 48,000 CIFAR10 images, finetune on 500 images with 5 poisons each time

