

Motivation

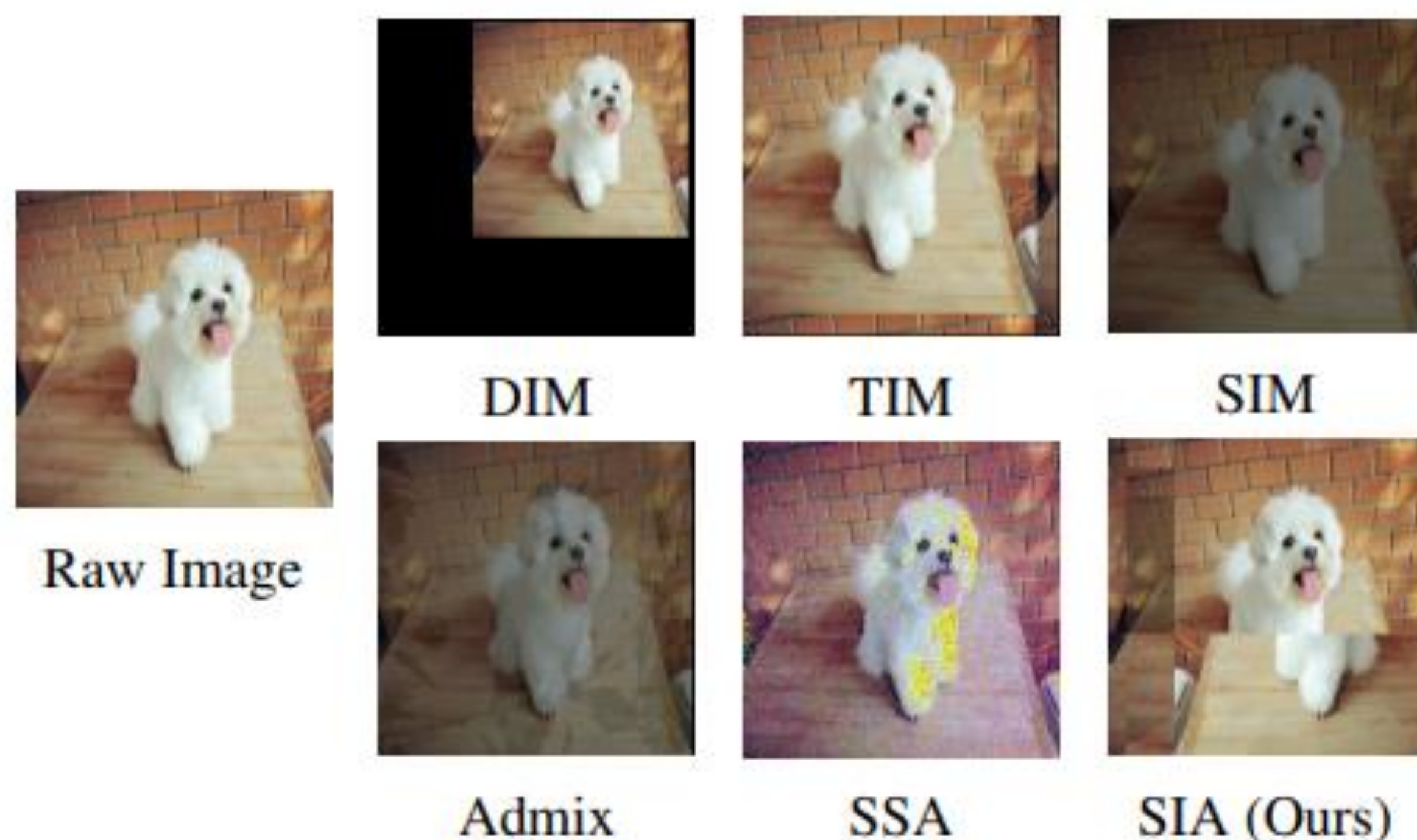
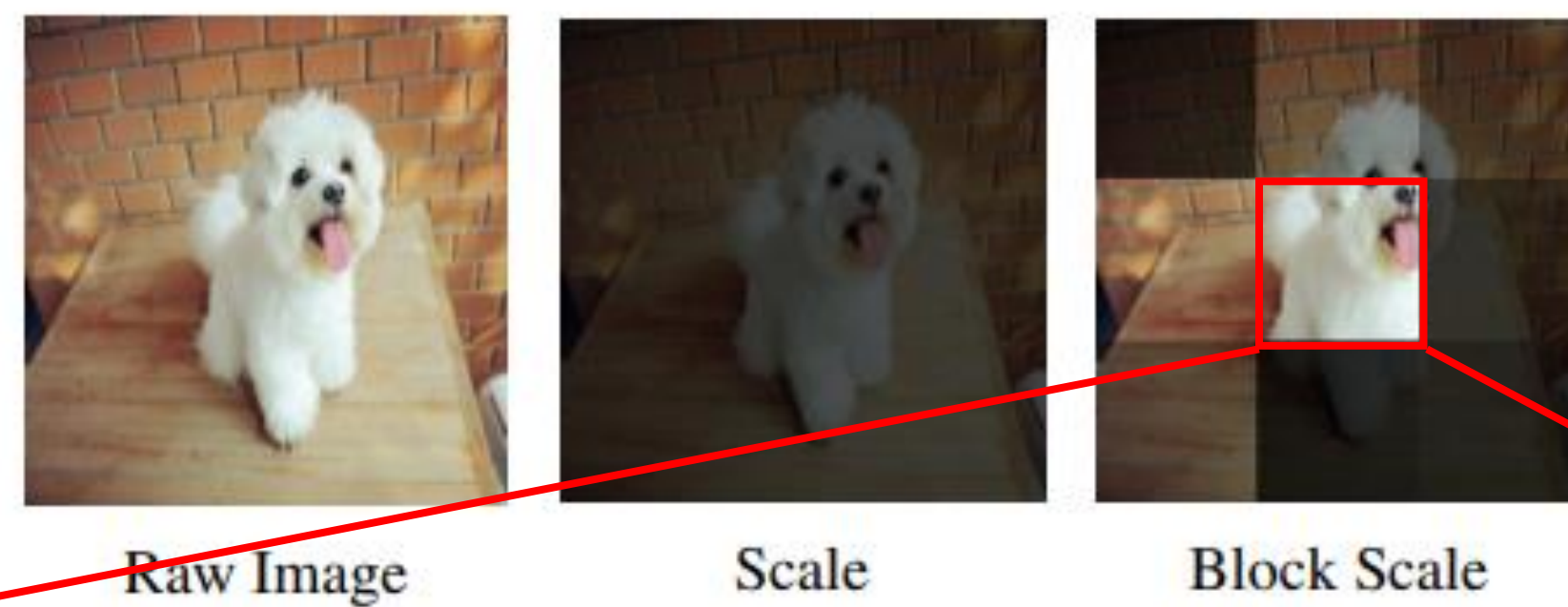
Assumption: Without harming the semantic information, the more diverse the transformed image is, the better transferability the adversarial examples have.

	TIM	DIM	SIM	SSA	Admix
Transferability	57.4	77.6	79.3	80.6	83.6
LPIPS	0.25	0.43	0.48	0.54	0.73

$$\text{LPIPS}(x, \hat{x}) = \frac{1}{H \times W} \sum_l \sum_{h,w} \|z_{h,w}^l - \hat{z}_{h,w}^l\|_2$$

The Structure of Image

Definition: Given an image x , which is randomly split into $s \times s$ blocks, the relative relation between each anchor point is the structure of image, where the anchor point is the center of the image block.



Structure Invariant Attack

Algorithm 1: Structure Invariant Attack

Input: Classifier $f(\cdot)$ with the loss function J ; The benign sample x with ground-truth label y ; The maximum perturbation ϵ , number of iterations T and decay factor μ ; Splitting number s ; Number of transformed images N

Output: An adversarial example.

- $\alpha = \epsilon/T$, $g_0 = 0$, $x_0^{adv} = x$
- for** $t = 0 \rightarrow T - 1$ **do**
- Constructing a set \mathcal{X} of N transformed images using SIT
- Calculating the average gradient on \mathcal{X} :

$$\bar{g}_{t+1} = \frac{1}{N} \sum_{x_i \in \mathcal{X}} \nabla_x J(x_i, y) \quad (2)$$

- Updating the momentum:

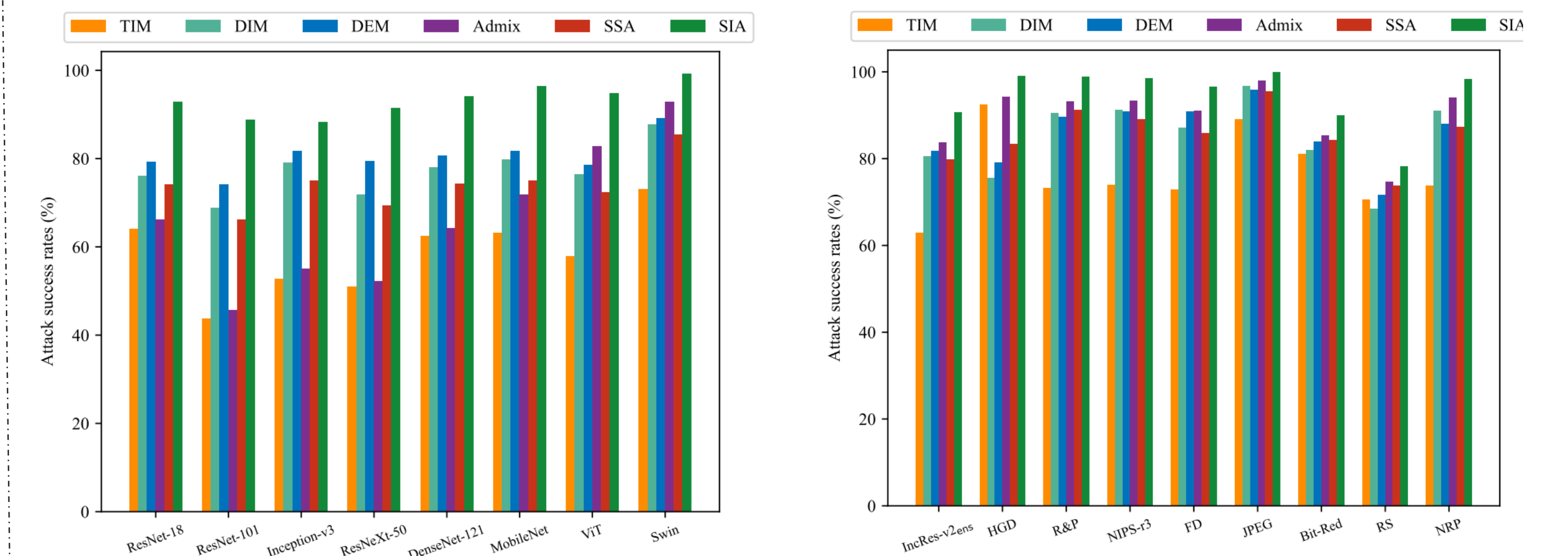
$$g_{t+1} = \mu g_t + \frac{\bar{g}_{t+1}}{\|\bar{g}_{t+1}\|_1} \quad (3)$$

- Updating the adversarial example:

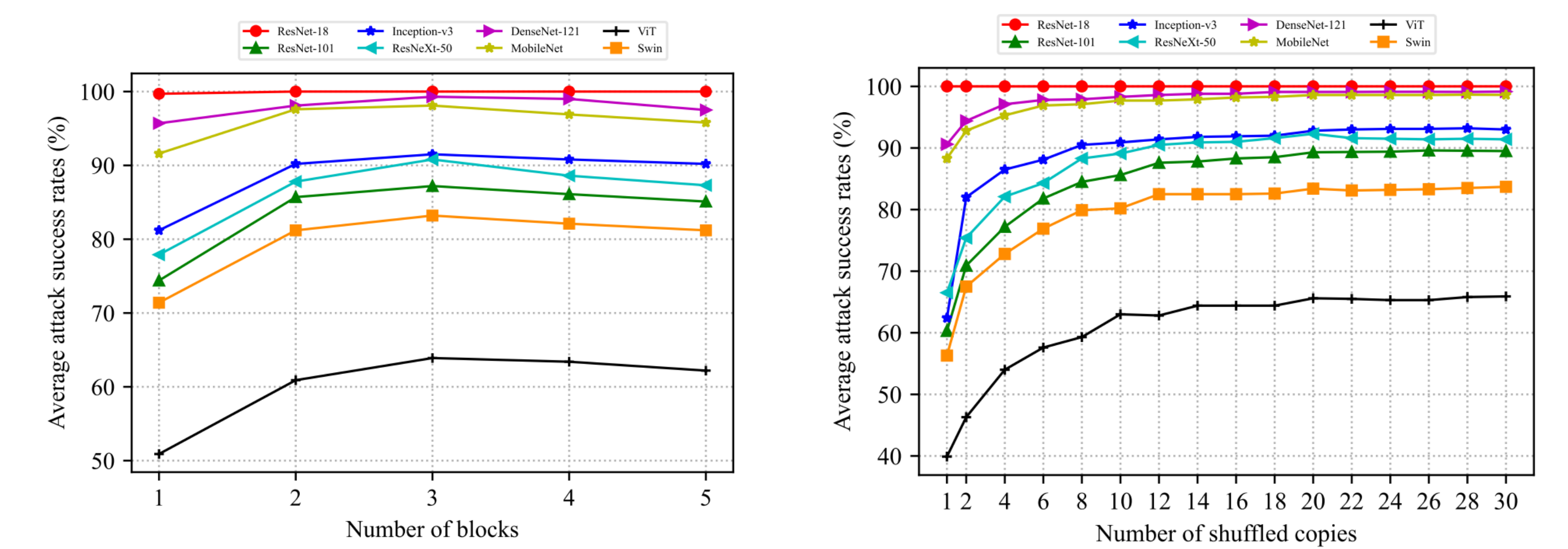
$$x_{t+1}^{adv} = \text{Clip}(x_t^{adv} + \alpha \cdot \text{sign}(g_{t+1}), 0, 1) \quad (4)$$

- return** x_T^{adv}

Attacking Ensemble Models and Defense Methods



Ablation study



	SIA	-VShift	-HShift	-VFlip	-HFlip	-Rotate	-Sclae	-Add Noise	-Resize	-DCT	-Dropout
92.1	89.7	90.1	90.1	88.4	88.3	90.1	90.6	90.2	90.1	90.7	

Attacking a Single Model

