

Introduction

- The **linear** (Goodfellow et al., ICLR 2015) and **non-flexible** (Fawzi et al., ICML 2015) nature of deep convolutional models makes them vulnerable to carefully crafted adversarial perturbations.
- Apart from attacking perspective, adversarial perturbations can be used to **measure** the linearity and flexibility of models, **regularize** models, and explore the **biases** of a model by analyzing the distances of samples to **decision boundaries**.
- RBF** networks have shown resilience against adversarial perturbations, but no successful **deep** RBF model has been trained yet.

Our hypothesis and proposed method

- Hypothesis:** A separable manifold should be resilient to perturbations which force a sample to cross the decision boundary.
- Proposed method:** Kernelized manifold transformation which leverages **RBF** to add **non-linearity** to models and learns a **transformation matrix** in Mahalanobis distance-like formulation to improve model **flexibility**.

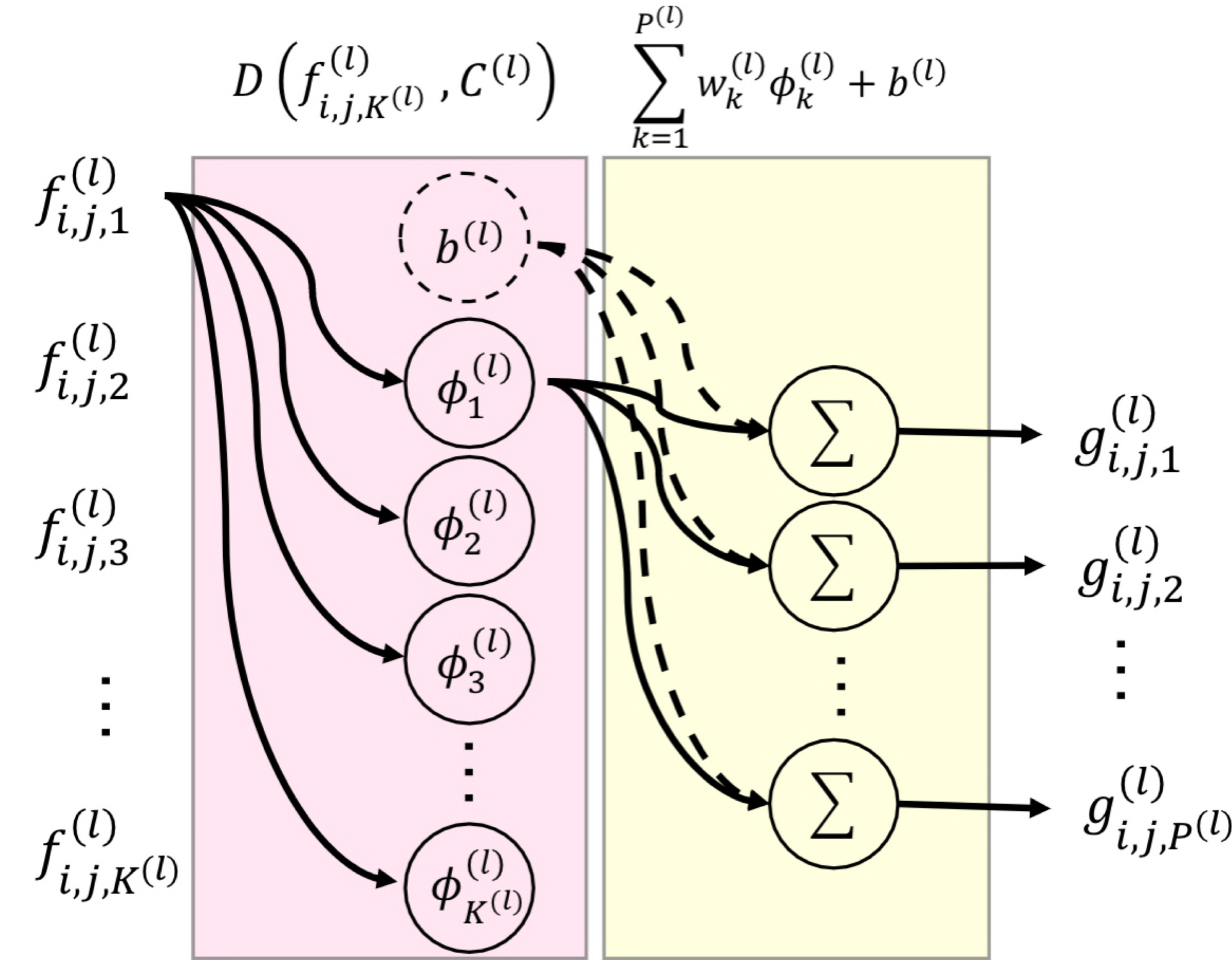
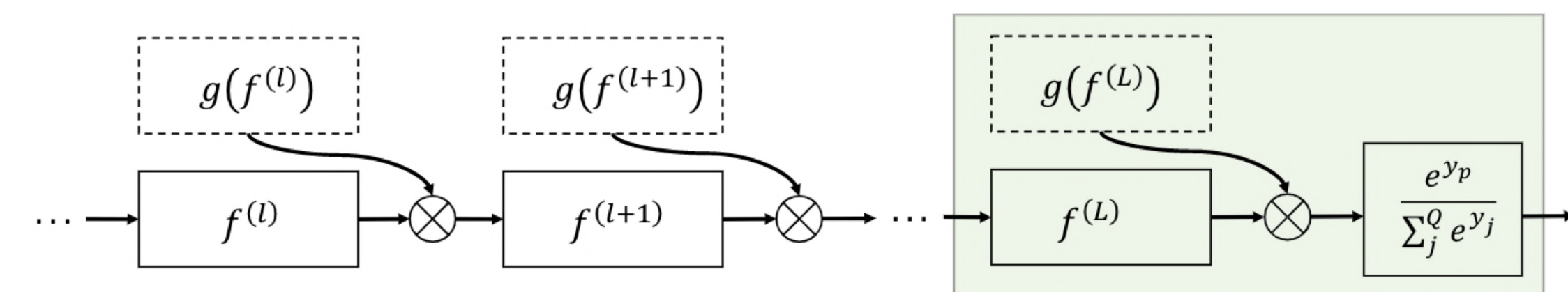
Input

$$f^{(l)} \in \mathcal{F}^{nm} \quad n \times m \times K^{(l)}$$

$$\Psi^{(l)} \in \mathbb{R}^{K^{(l)} \times K^{(l)}}, \quad \kappa : \mathcal{F} \times \mathcal{F} \rightarrow \mathbb{R}, \quad \beta, \quad C, \quad W$$

Output

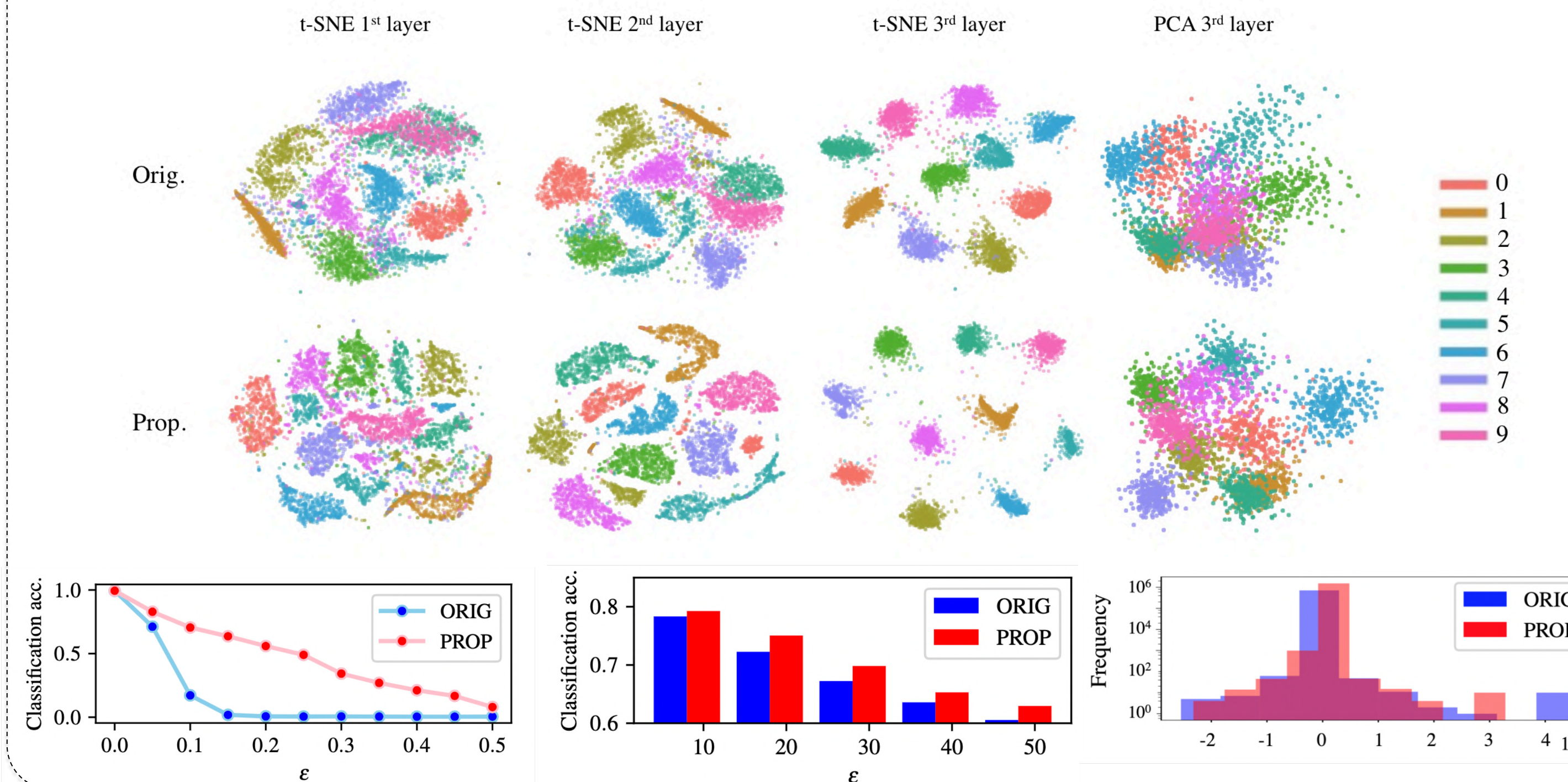
$$g^{(l)} \in \mathcal{G}^{nm} \quad n \times m \times P^{(l)}$$



$$\phi_k^{(l)} = e^{-\beta_k^{(l)} D(f_{K^{(l)}}^{(l)}, c_k^{(l)})}$$

$$D(f_{K^{(l)}}^{(l)}, c_k^{(l)}) = (f_{K^{(l)}}^{(l)} - c_k^{(l)})^T (\Psi^{(l)})^{-1} (f_{K^{(l)}}^{(l)} - c_k^{(l)})$$

Results

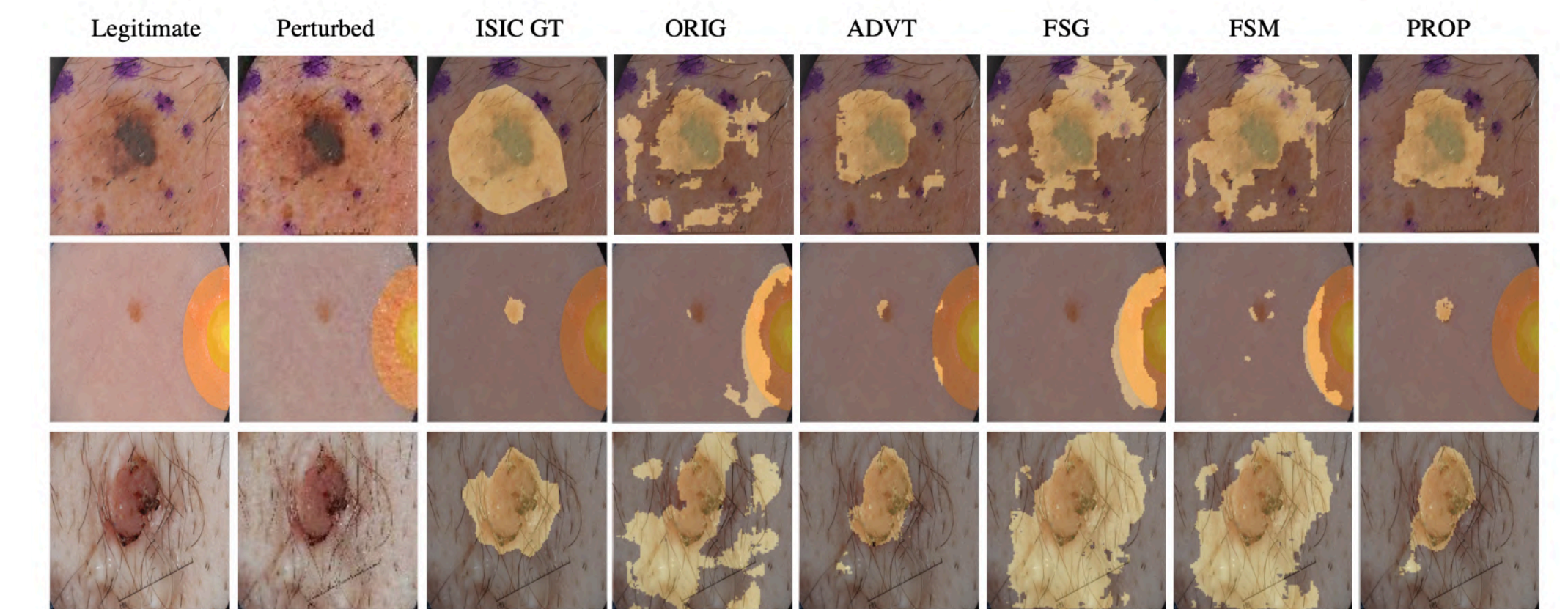


Classification (MNIST)

Models	Clean	L_2		L_∞				
		C&W [5]	GN [33]	FGSM [19]	BIM [19]	MIM [9]	PGD [27]	SPSA [49]
ORIG [30]	0.9930	0.1808	0.7227	0.0968	0.0070	0.0051	0.1365	0.3200
Binary CNN [23]	0.9850	n/a	0.9200	0.7100	0.7000	0.7000	n/a	n/a
NN [23]	0.9690	n/a	0.9100	0.6800	0.4300	0.2600	n/a	n/a
Binary ABS [23]	0.9900	n/a	0.8900	0.8500	0.8600	0.8500	n/a	n/a
ABS [23]	0.9900	n/a	0.9800	0.3400	0.1300	0.1700	n/a	n/a
Fortified Net [20]	0.9893	0.6058	n/a	0.9131	n/a	n/a	0.7954	n/a
PROP	0.9942	0.9879	0.7506	0.8582	0.7887	0.6425	0.8157	0.7092

White-box (segmentation)

Network	Method	Clean	10i (% Accuracy drop)	30i (% Accuracy drop)
U-Net [34]	ORIG [34]	0.7743 ± 0.0202	0.5594 ± 0.0196(27.75%)	0.4396 ± 0.0222(43.23%)
	FSG [55]	0.7292 ± 0.0229	0.6382 ± 0.0206(15.58%)	0.5858 ± 0.0218(24.34%)
	FSM [55]	0.7695 ± 0.0198	0.6039 ± 0.0199(22.01%)	0.5396 ± 0.0211(30.31%)
	ADVT [14]	0.6703 ± 0.0273	0.7012 ± 0.0255(9.44%)	0.6700 ± 0.0260(13.47%)
	PROP	0.7780 ± 0.0209	0.7619 ± 0.0208 (1.60%)	0.7248 ± 0.0226 (6.39%)
V-Net [29]	ORIG [34]	0.8070 ± 0.0189	0.5320 ± 0.0207(34.10%)	0.3865 ± 0.0217(52.10%)
	FSG [55]	0.7886 ± 0.0205	0.6990 ± 0.0189(13.38%)	0.6840 ± 0.0188(15.24%)
	FSM [55]	0.8084 ± 0.0189	0.5928 ± 0.0209(26.54%)	0.5144 ± 0.0218(36.26%)
	ADVT [14]	0.7924 ± 0.0162	0.7121 ± 0.0174(11.76%)	0.7113 ± 0.0179 (11.85%)
	PROP	0.8213 ± 0.0177	0.7384 ± 0.0169 (8.50%)	0.6944 ± 0.0178(13.95%)



Black-box (segmentation)

	U-Net [34]	U-PROP	V-Net [29]	V-PROP
-	-	-	-	-
U-Net [34]	-	0.7341 ± 0.0205	0.6364 ± 0.0189	0.7210 ± 0.0189
U-PROP	0.7284 ± 0.0219	-	0.6590 ± 0.0218	0.7262 ± 0.0241
V-Net [29]	0.7649 ± 0.0168	0.7773 ± 0.0167	-	0.7478 ± 0.2090
V-PROP	0.7922 ± 0.0188	0.7964 ± 0.0192	0.6948 ± 0.0171	-

Classification (chest x-ray)

		Defense				
Attack	Iteration	ORIG	GDA	FSM	PROP	
L_1 BIM [19]	5	0	0	0.55	0.63	
L_∞ BIM [19]	5	0	0	0.54	0.65	
Clean	-	0.74	0.75	0.57	0.74	