

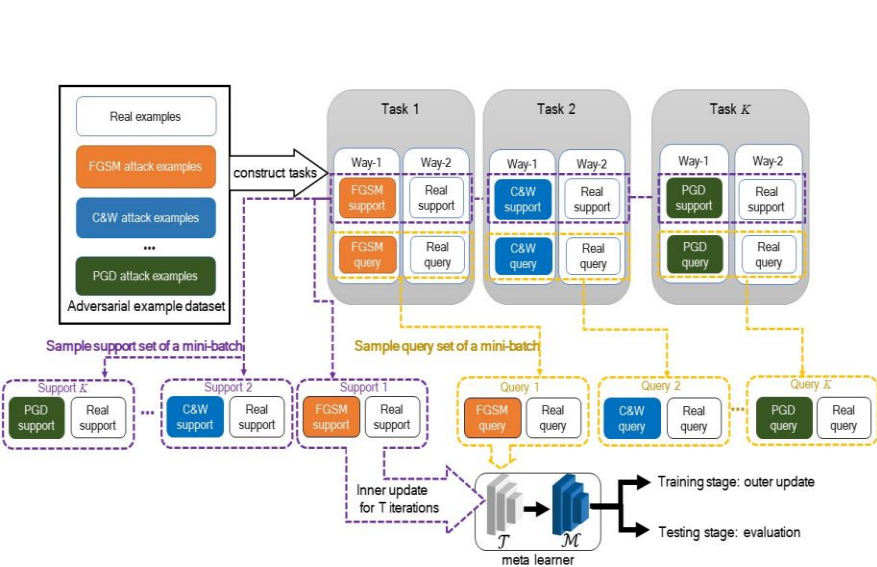
# MetaAdvDet: Towards Robust Detection of Evolving Adversarial Attacks

**Chen Ma**<sup>1</sup>, Chenxu Zhao<sup>2</sup>, Hailin Shi<sup>2</sup>, Li Chen<sup>1</sup>, Junhai Yong<sup>1</sup>, Dan Zeng<sup>3</sup>

1. School of Software, Tsinghua University, Beijing, China

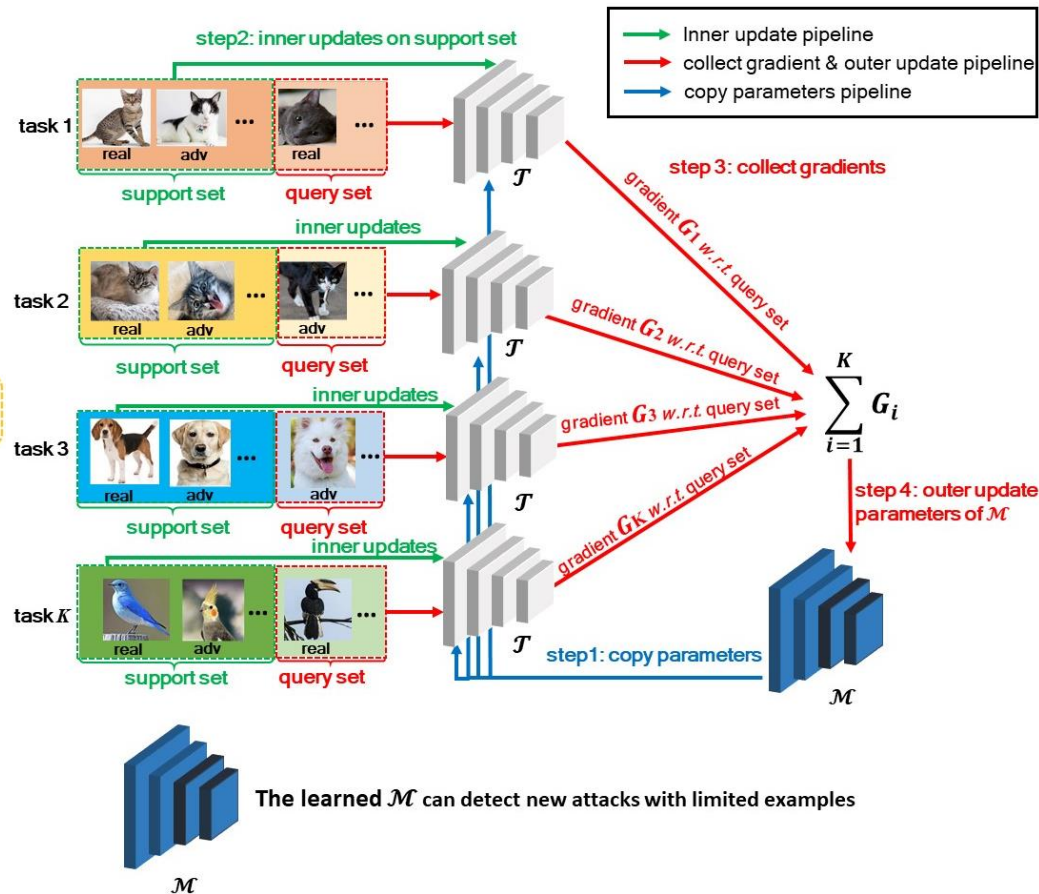
2. JD AI Research

3. Shanghai University



MetaAdvDet uses the meta-learning based method with a double-network framework to learn the capability of detecting the evolving adversarial attacks.

Paper ID: **P2B-05**



# Results

Table 1: cross-adversary benchmark

Dataset	Method	F1 score	
		1-shot	5-shot
AdvCIFAR	DNN	0.495	0.639
	DNN (balanced)	0.536	0.643
	NeuralFP [8]	<b>0.698</b>	0.700
	TransformDet [45]	0.662	0.697
	MetaAdvDet (ours)	0.685	<b>0.791</b>
AdvMNIST	DNN	0.812	0.852
	DNN (balanced)	0.797	0.808
	NeuralFP [8]	0.780	0.906
	TransformDet [45]	0.840	0.904
	MetaAdvDet (ours)	<b>0.987</b>	<b>0.993</b>
AdvFashionMNIST	DNN	0.782	0.885
	DNN (balanced)	0.744	0.850
	NeuralFP [8]	0.798	0.817
	TransformDet [45]	0.712	0.879
	MetaAdvDet (ours)	<b>0.848</b>	<b>0.944</b>

Paper ID: **P2B-05**

Table 2: cross-domain benchmark

Train Domain	Test Domain	Method	F1 score	
			1-shot	5-shot
AdvMNIST	AdvFashionMNIST	DNN (balanced)	0.698	0.813
		NeuralFP [8]	0.748	0.811
		TransformDet [45]	0.664	0.808
		MetaAdvDet (ours)	<b>0.799</b>	<b>0.870</b>
AdvFashionMNIST	AdvMNIST	DNN (balanced)	0.950	0.977
		NeuralFP [8]	0.775	0.836
		TransformDet [45]	0.934	0.940
		MetaAdvDet (ours)	<b>0.956</b>	<b>0.981</b>

Table 3: white-box benchmark

Dataset	Method	I-FGSM Attack		C&W Attack	
		1-shot	5-shot	1-shot	5-shot
CIFAR-10	DNN (balanced)	0.466	0.537	0.459	0.527
	TransformDet [45]	<b>0.593</b>	<b>0.728</b>	0.443	0.502
	MetaAdvDet (ours)	0.553	0.633	<b>0.548</b>	<b>0.607</b>
MNIST	DNN (balanced)	0.857	0.956	0.814	0.913
	TransformDet [45]	0.864	0.952	0.775	0.893
	MetaAdvDet (ours)	<b>0.968</b>	<b>0.994</b>	<b>0.920</b>	<b>0.990</b>
FashionMNIST	DNN (balanced)	0.745	0.890	0.726	0.853
	TransformDet [45]	0.837	0.920	0.747	0.853
	MetaAdvDet (ours)	<b>0.849</b>	<b>0.963</b>	<b>0.882</b>	<b>0.967</b>

# Thanks!