

# 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>

<sup>1</sup> School of Software, Tsinghua University, Beijing, China

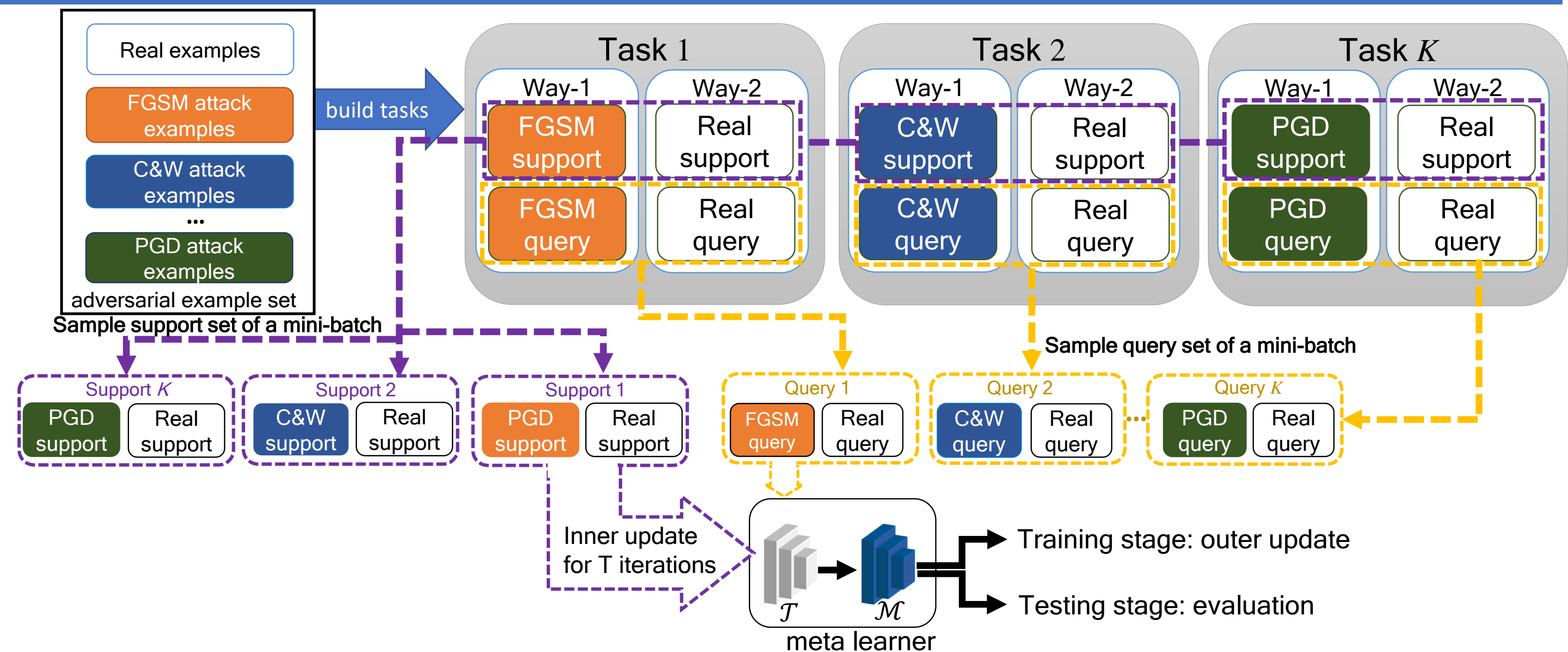
<sup>2</sup> JD AI research      <sup>3</sup> Shanghai University

## Motivations

The shortcomings of existing approaches for detecting the evolving adversarial attacks:

- ◆ Labeled samples of new attacks are insufficient and expensive.
- ◆ New attacks evolve much faster than the high-cost data collection.
- ◆ The small scale of data leads to the few-shot learning problem → MetaAdvDet is proposed.

## Learning from the few-shot tasks



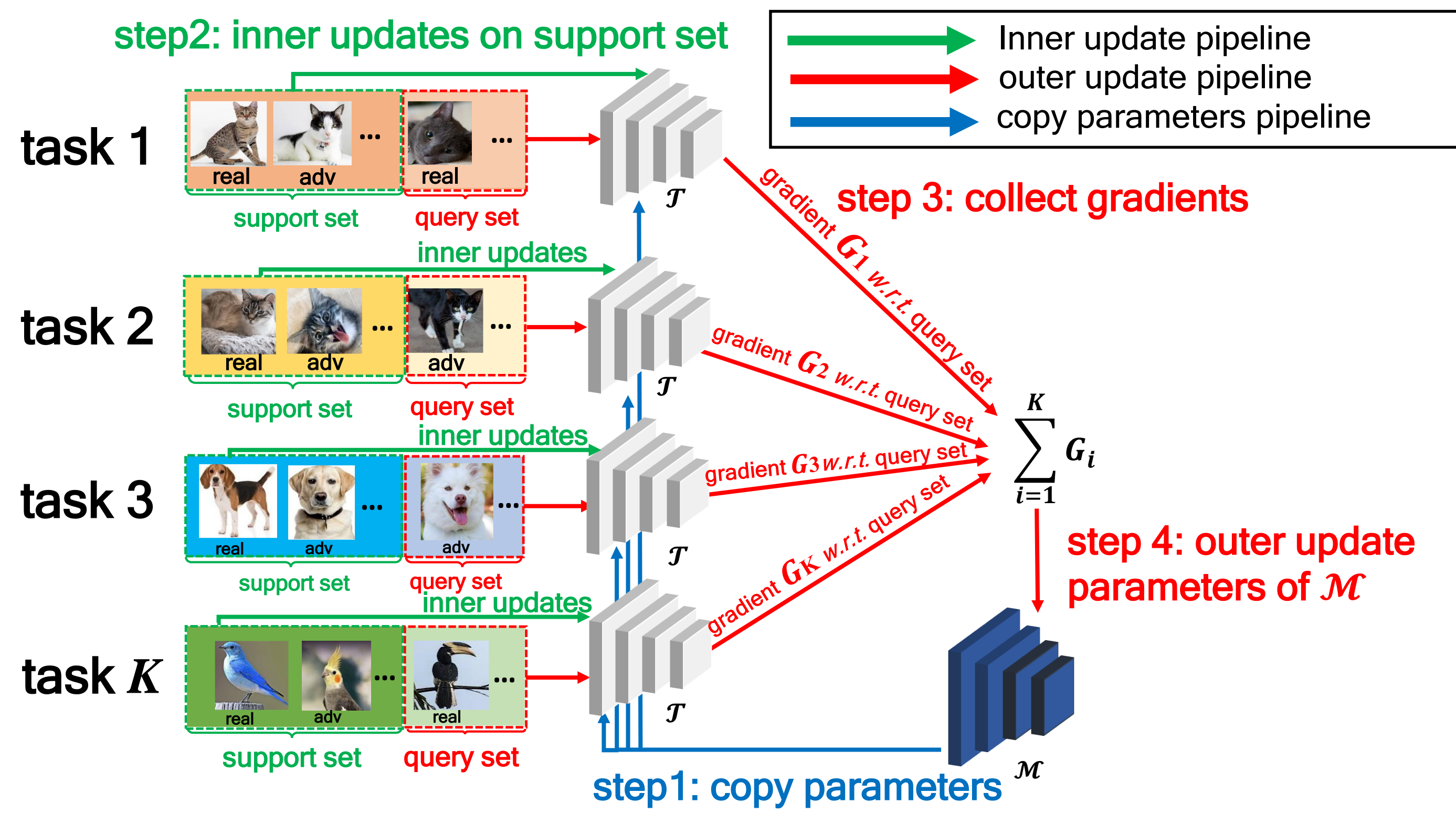
## Proposed Benchmark

Benchmark	Test Protocols	
Datasets	CIFAR-10, MNIST and FashionMNIST	
Cross-Adversary Benchmark (simulate the situation of evolving attacks)	Train Adversary	Test Adversary
	FGSM, MI-FGSM, BIM, PGD, C&W, JSMA, SPSA, VAT, MaxConfidence	EAD, semantic, DeepFool, Spatial Transformation, NewtonFool
Cross-Domain Benchmark	Train Domain	Test Domain
	MNIST FashionMNIST	FashionMNIST MNIST
Cross-Architecture Benchmark (evaluate the detection of adversarial examples with new architecture)	Train Architecture	Test Architecture
	ResNet-10 ResNet-18 Conv-4 ResNet-10	ResNet-18 ResNet-10 ResNet-10 Conv-4

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>

## Results of cross-domain benchmark

## Method



The learned  $\mathcal{M}$  can detect new attacks with limited examples

MetaAdvDet is equipped with a double-network framework  $\mathcal{M}$  and  $\mathcal{T}$ .  $\mathcal{T}$  focuses on learning from individual tasks. After a couple of iterations,  $\mathcal{T}$  converges and computes the gradient  $G_i$  which are accumulated by  $\mathcal{M}$  to update  $\mathcal{M}$ 's parameters for achieving the fast adaption capability in detecting new attacks.

## Results

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>

## Results of cross-adversary 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>

## Results of white-box attack benchmark