



Machine learning security

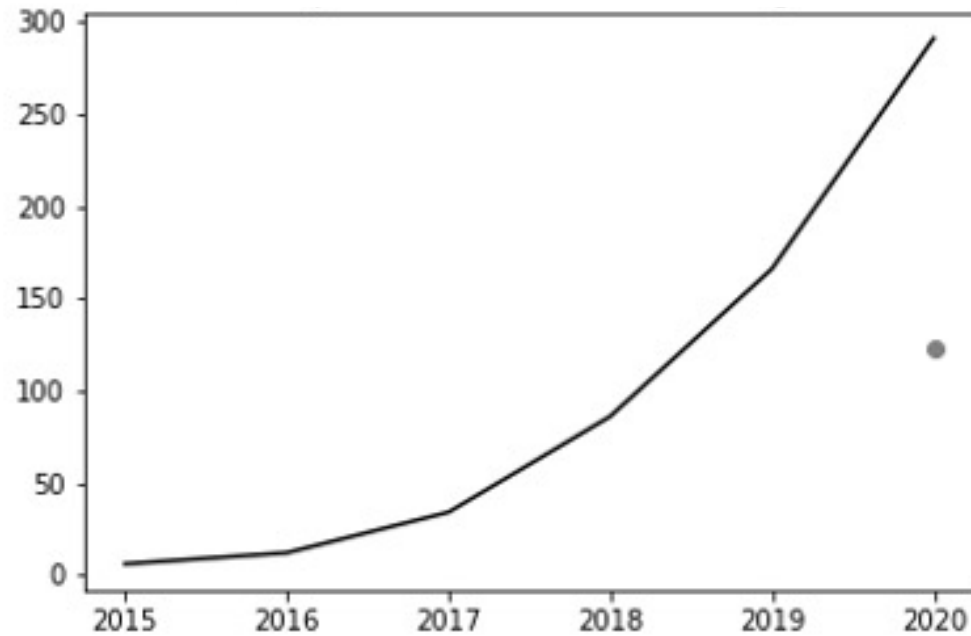
Danila Yu. Emelyanov

Lead mathematician, JSC “NPK Kryptonite”

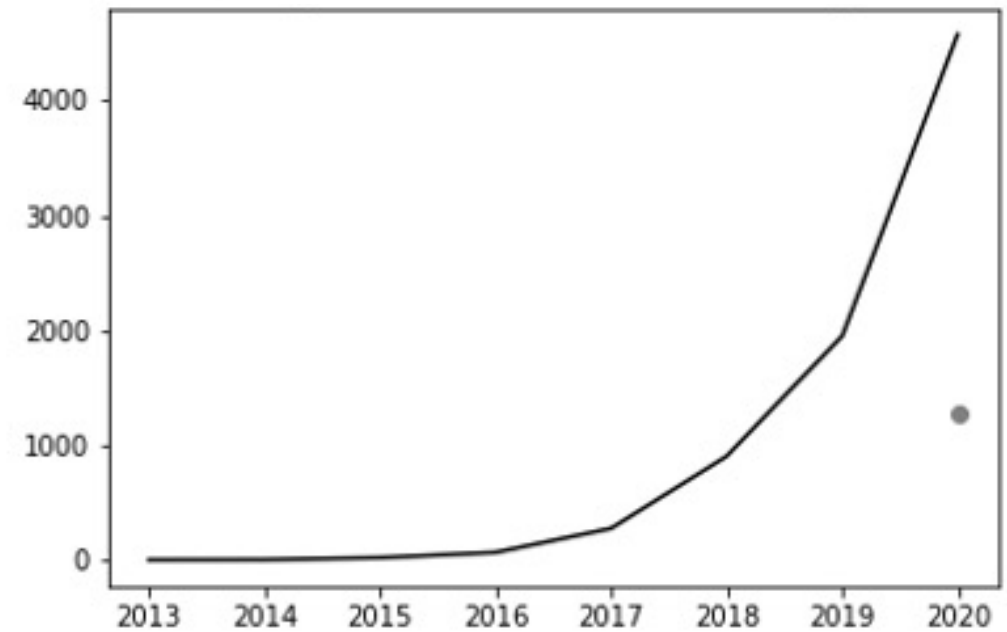


Different search requests results count on arxiv.org

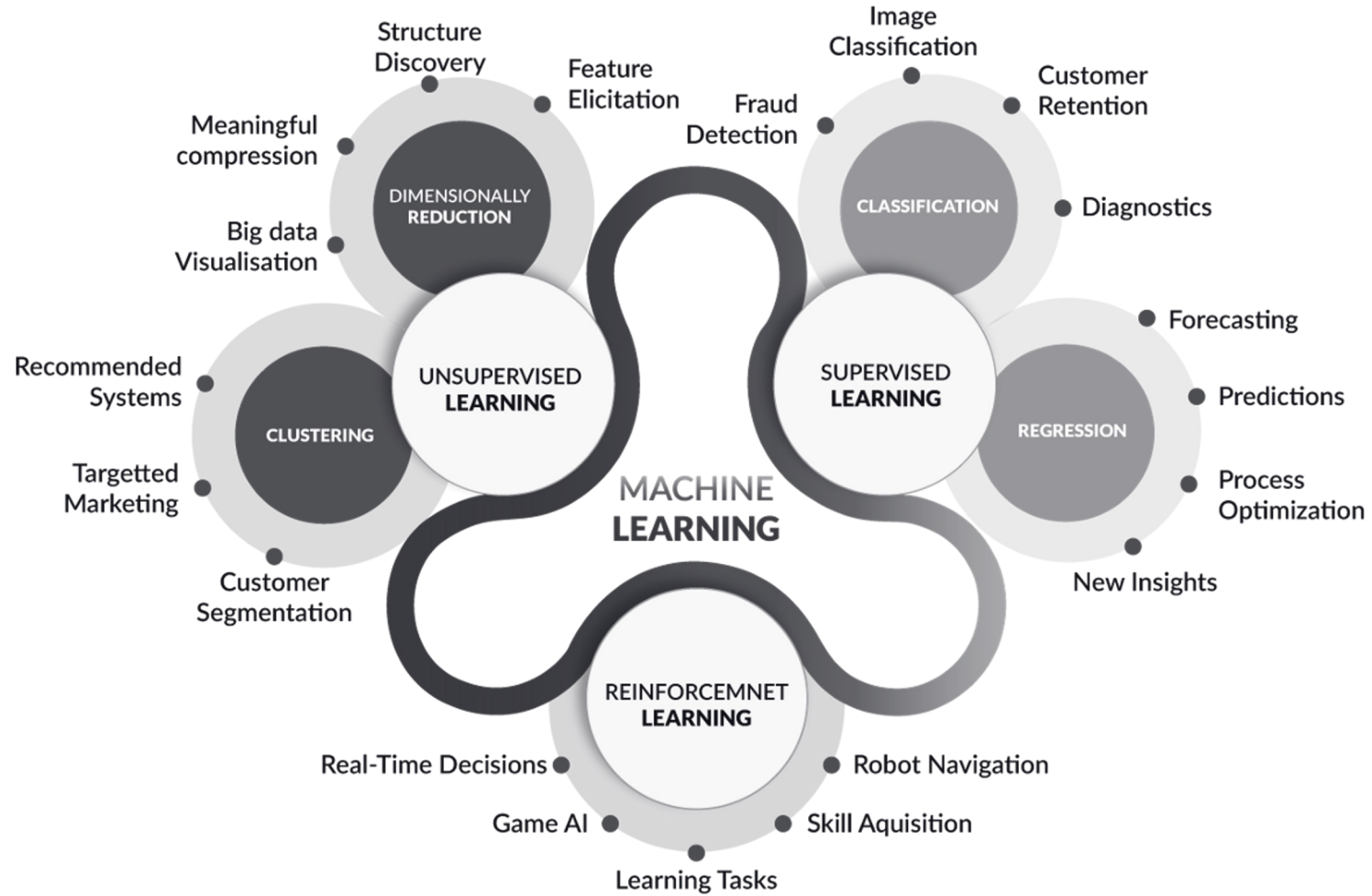
"threat machine learning"



"adversarial machine learning"

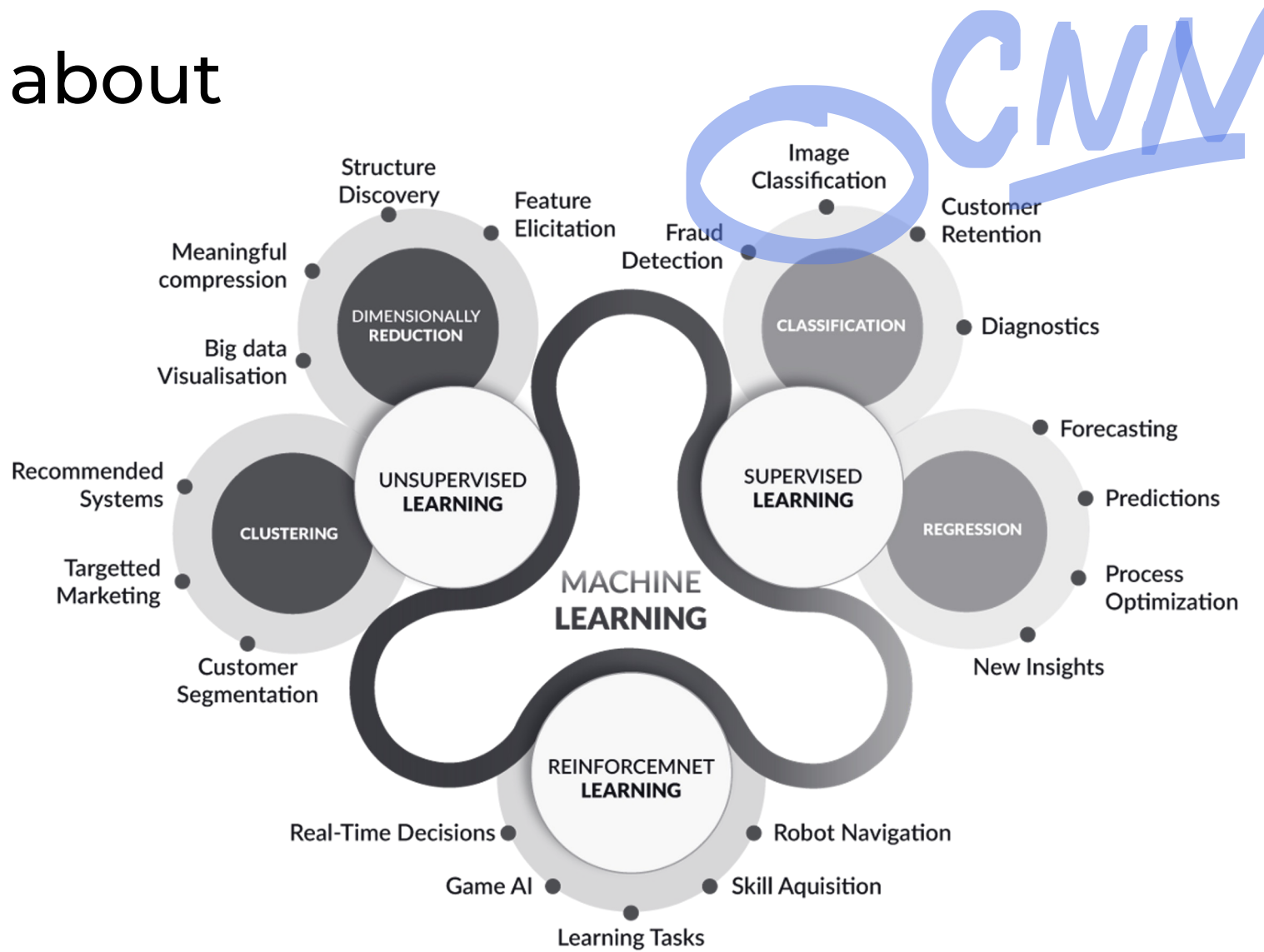


Let's talk about

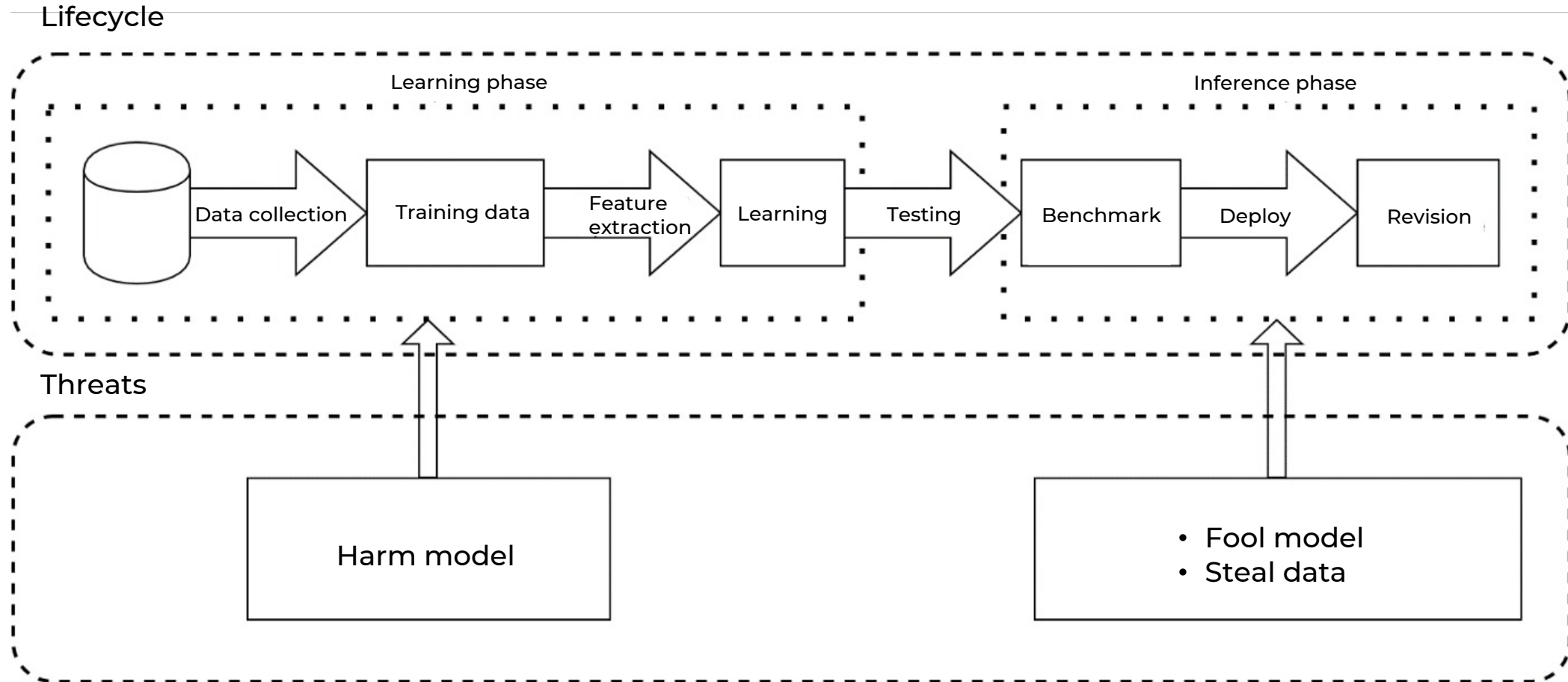


Let's talk about

CNN-based supervised learning model for image classification

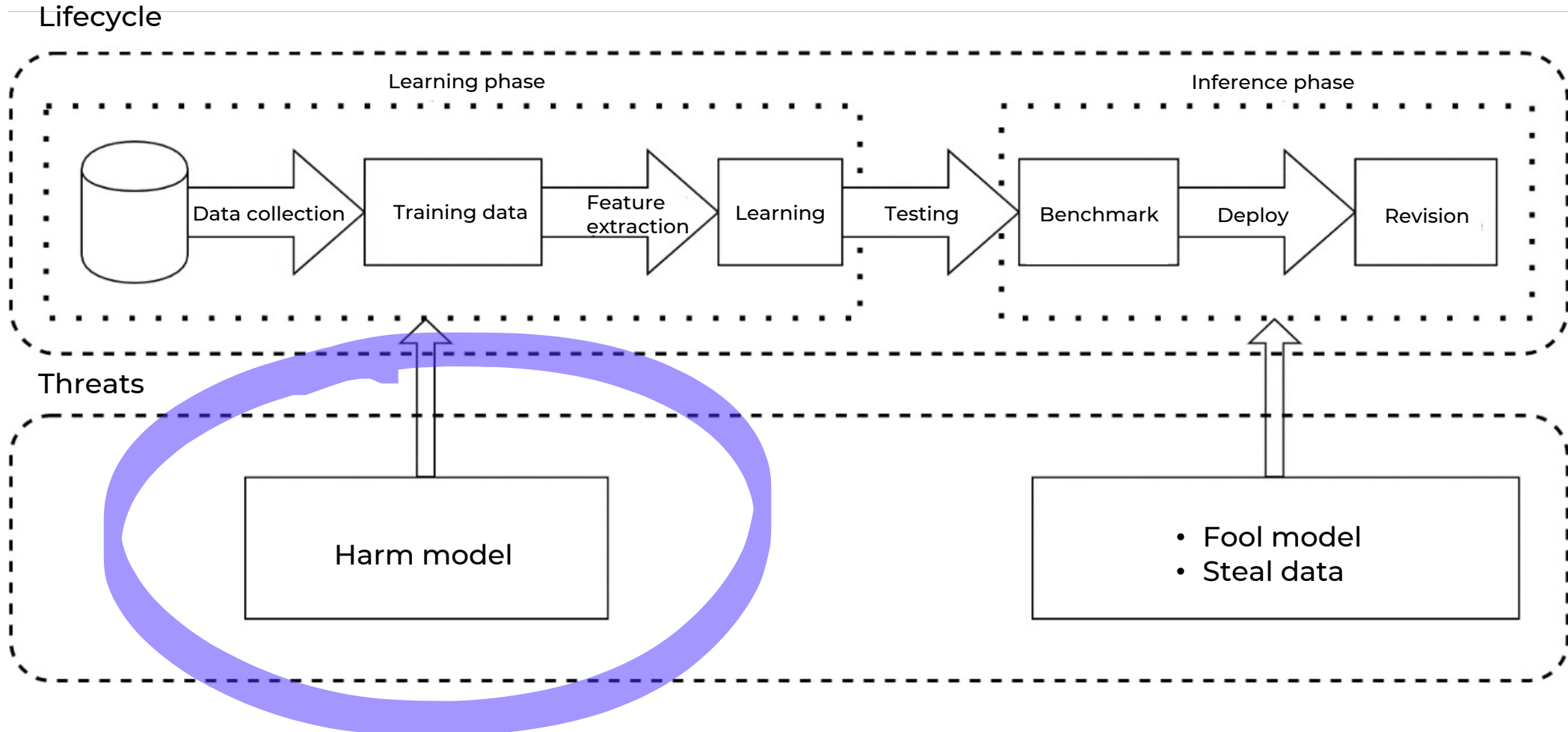


Threat model



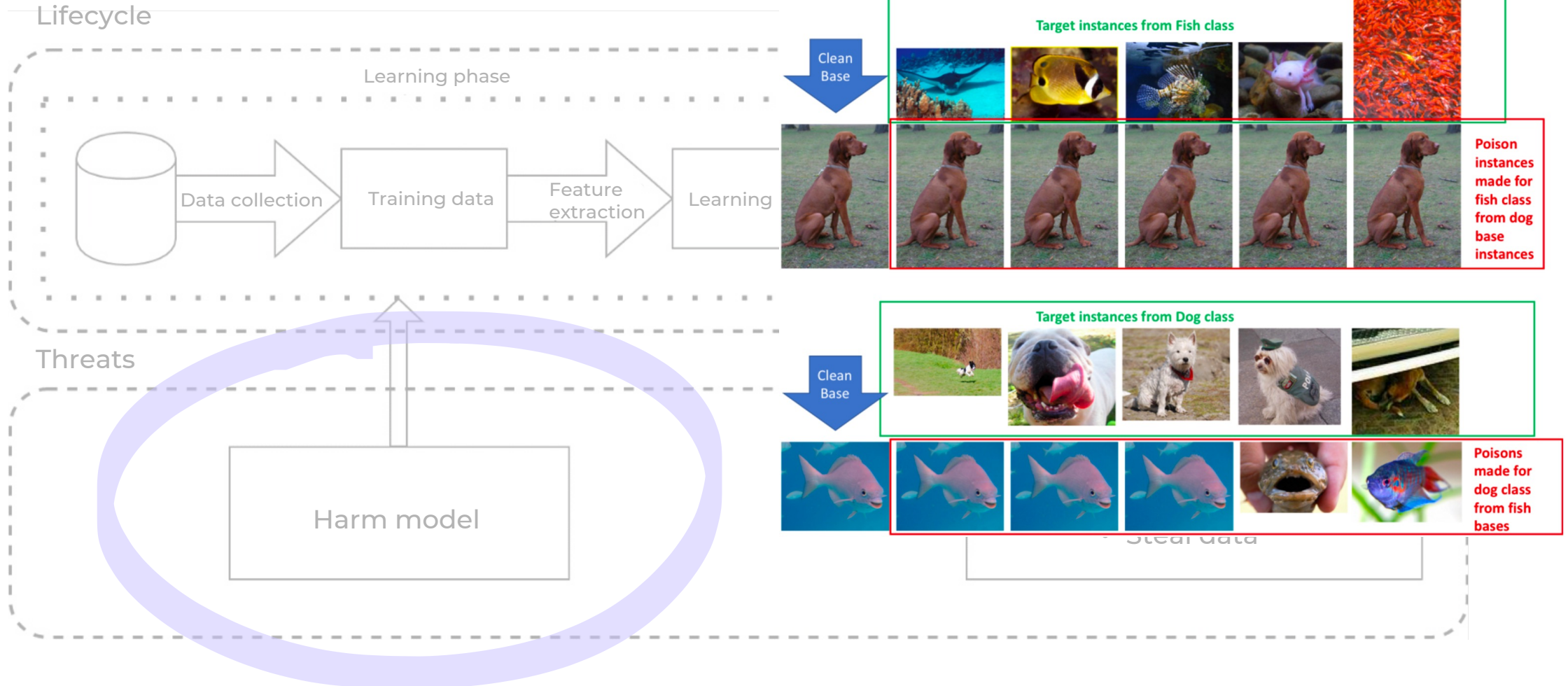
Threat model

Training phase



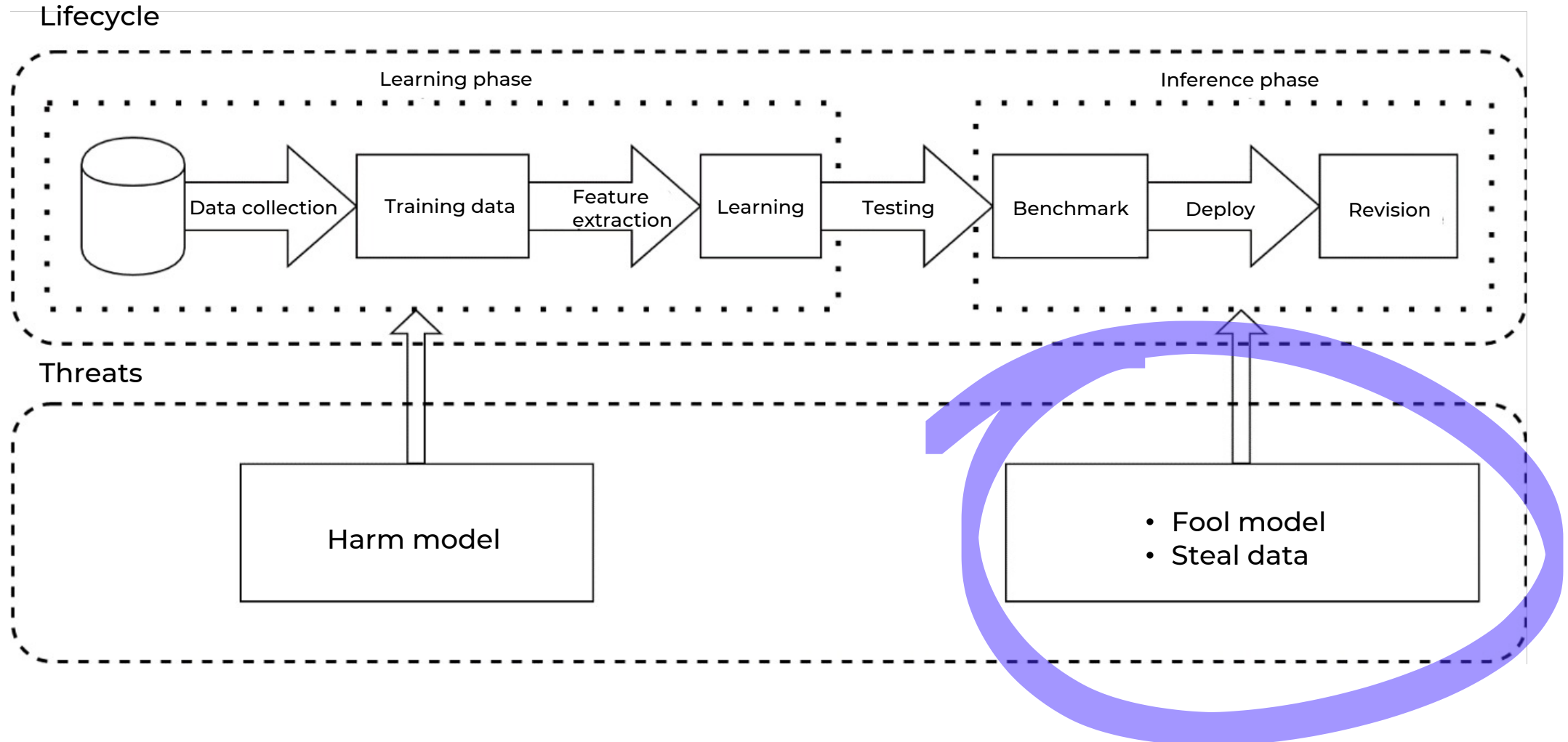
Threat model

Training phase



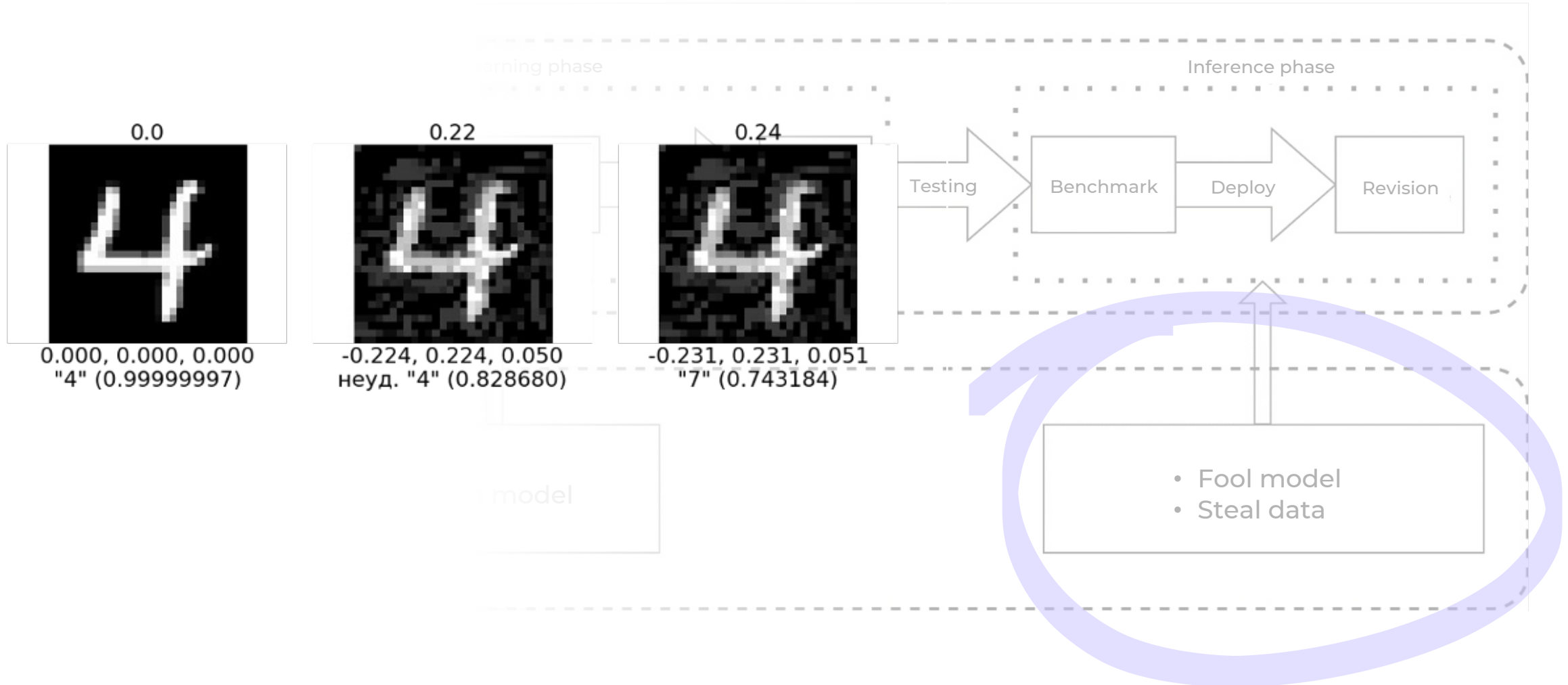
Threat model

Inference phase



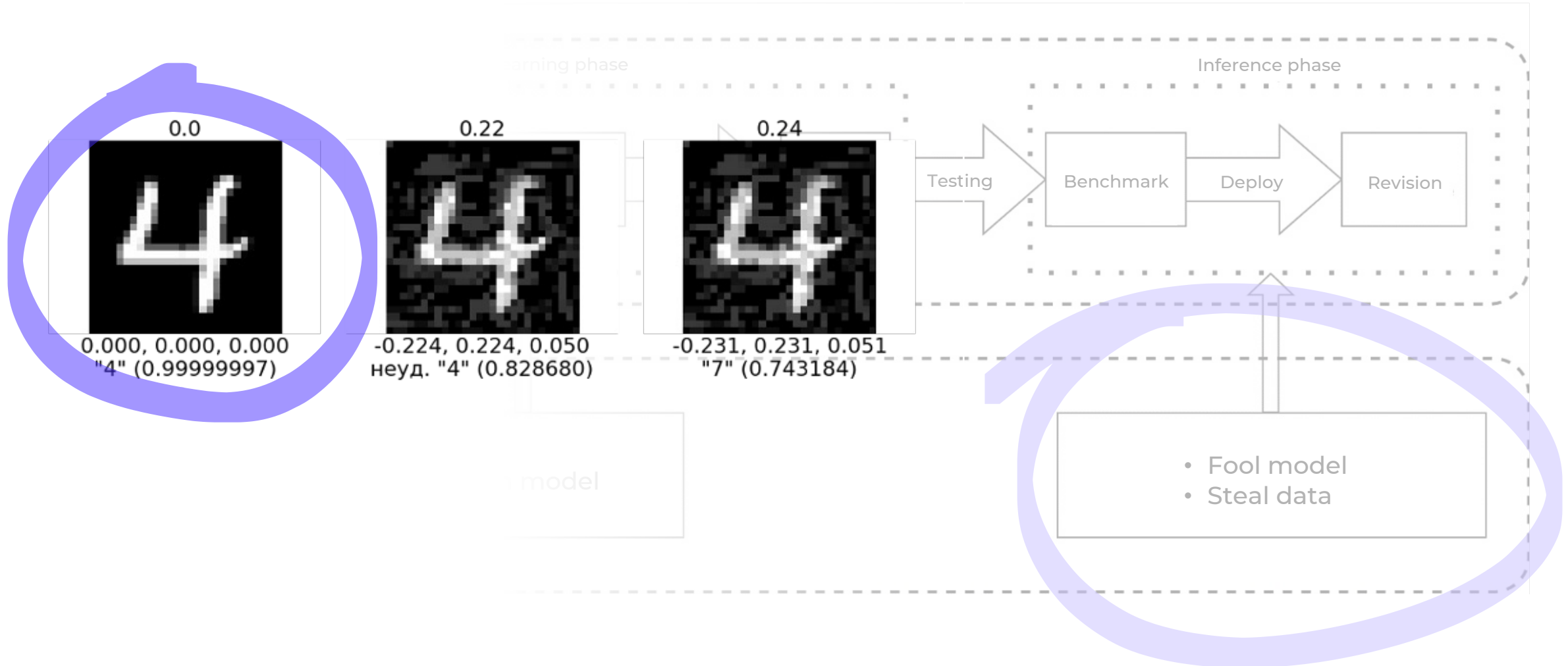
Adversarial examples

Perturbations



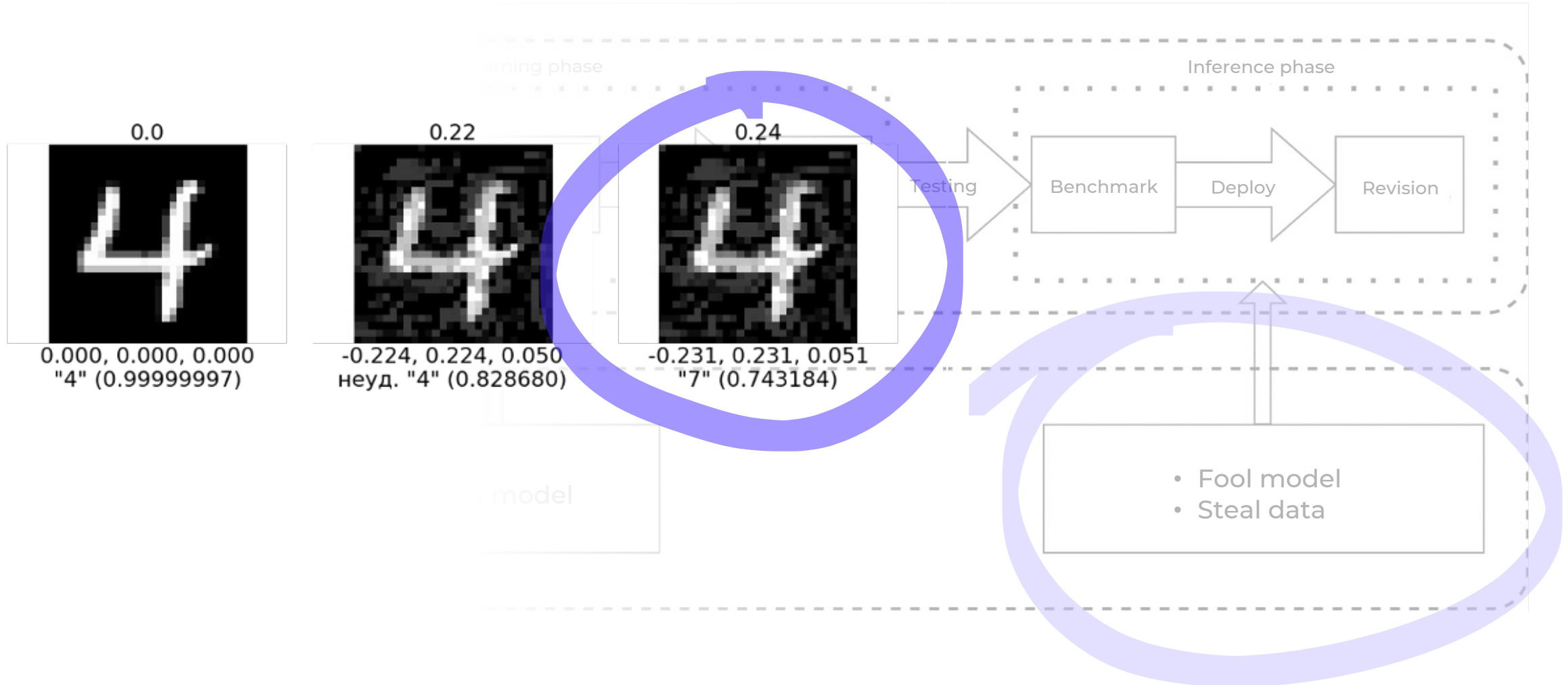
Adversarial examples

Perturbations



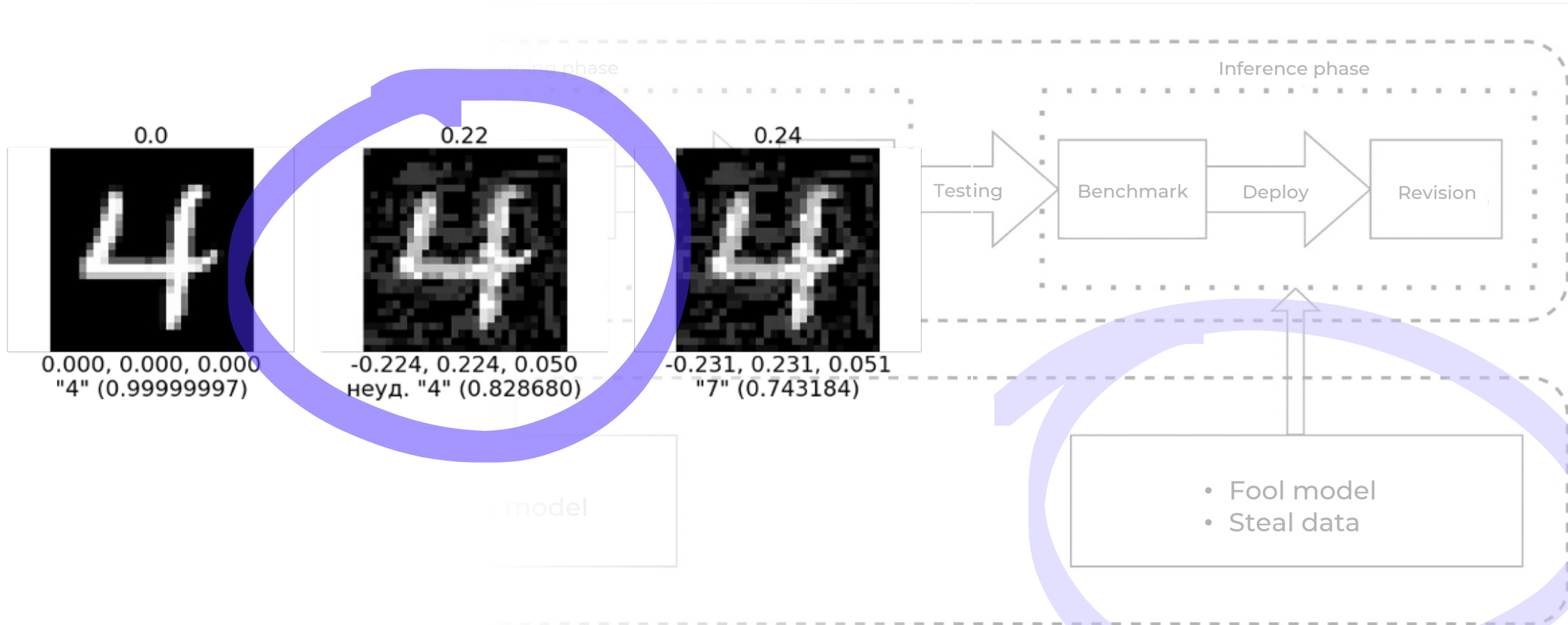
Adversarial examples

Perturbations



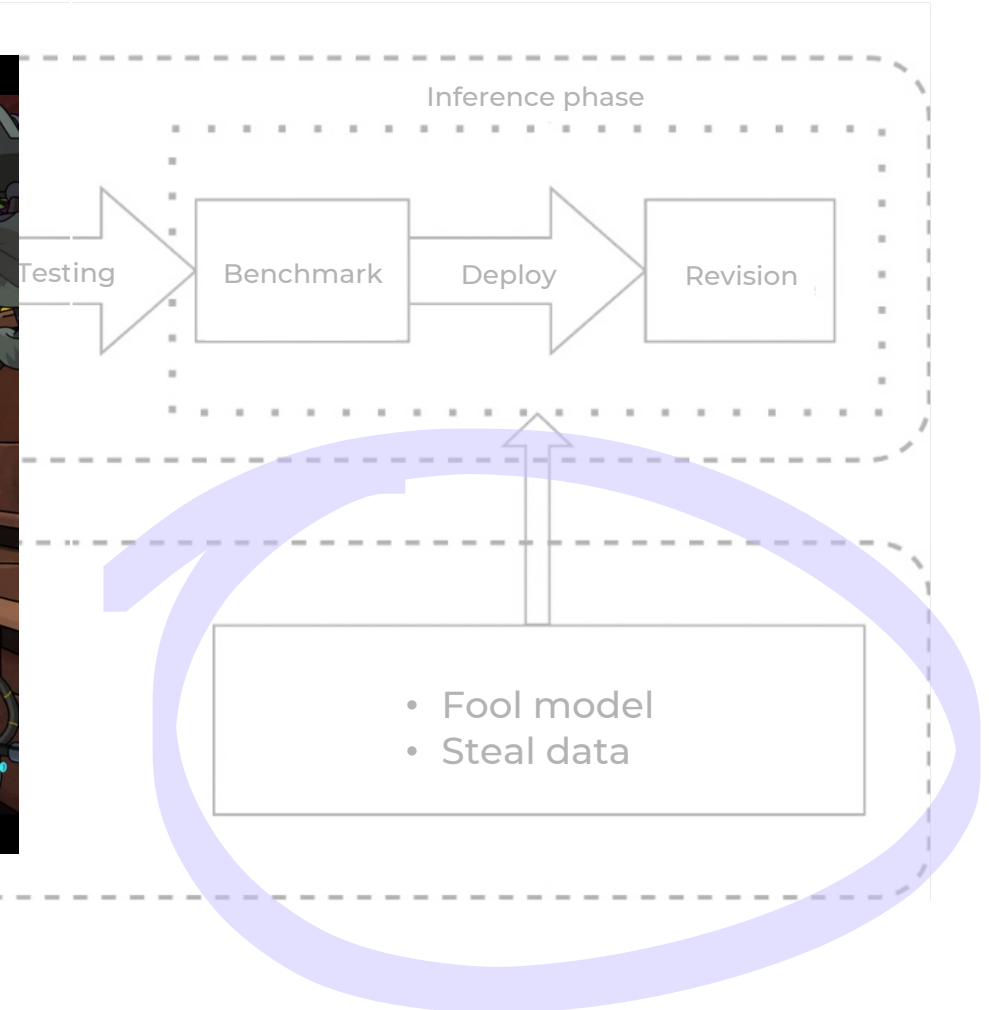
Adversarial examples

Perturbations



Adversarial examples

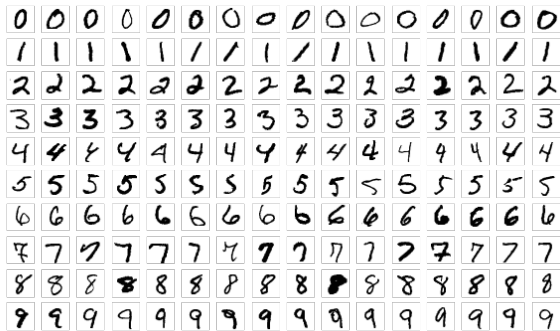
Watermarks



Summary

We shall deal with

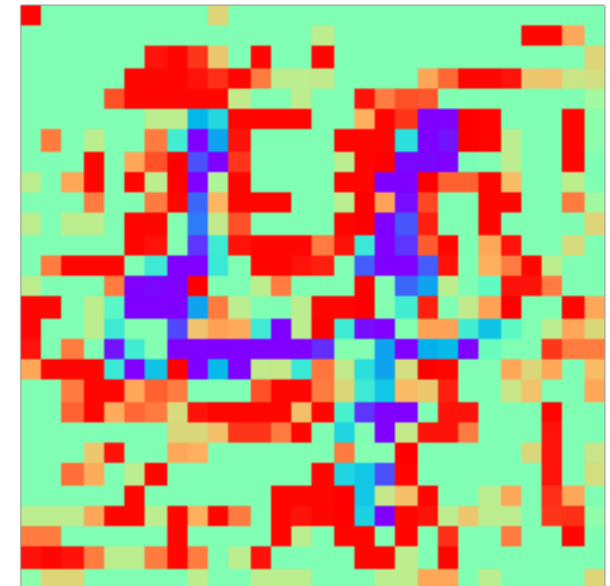
MNIST, PyTorch



Adversarial examples \equiv perturbations



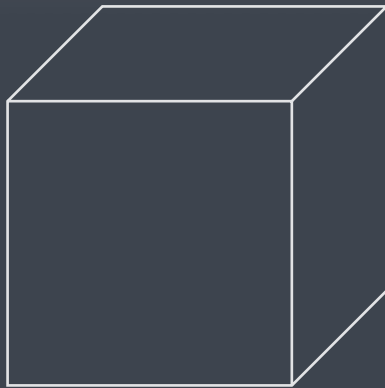
-0.231, 0.231, 0.051
"7" (0.743184)



-0.231, 0.231, 0.051
"7" (0.743184)

Black box, gray box, white box

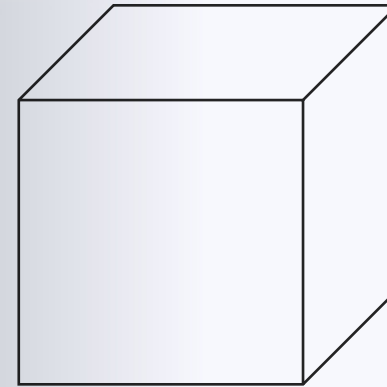
Adversary knowledge about the model



Nothing

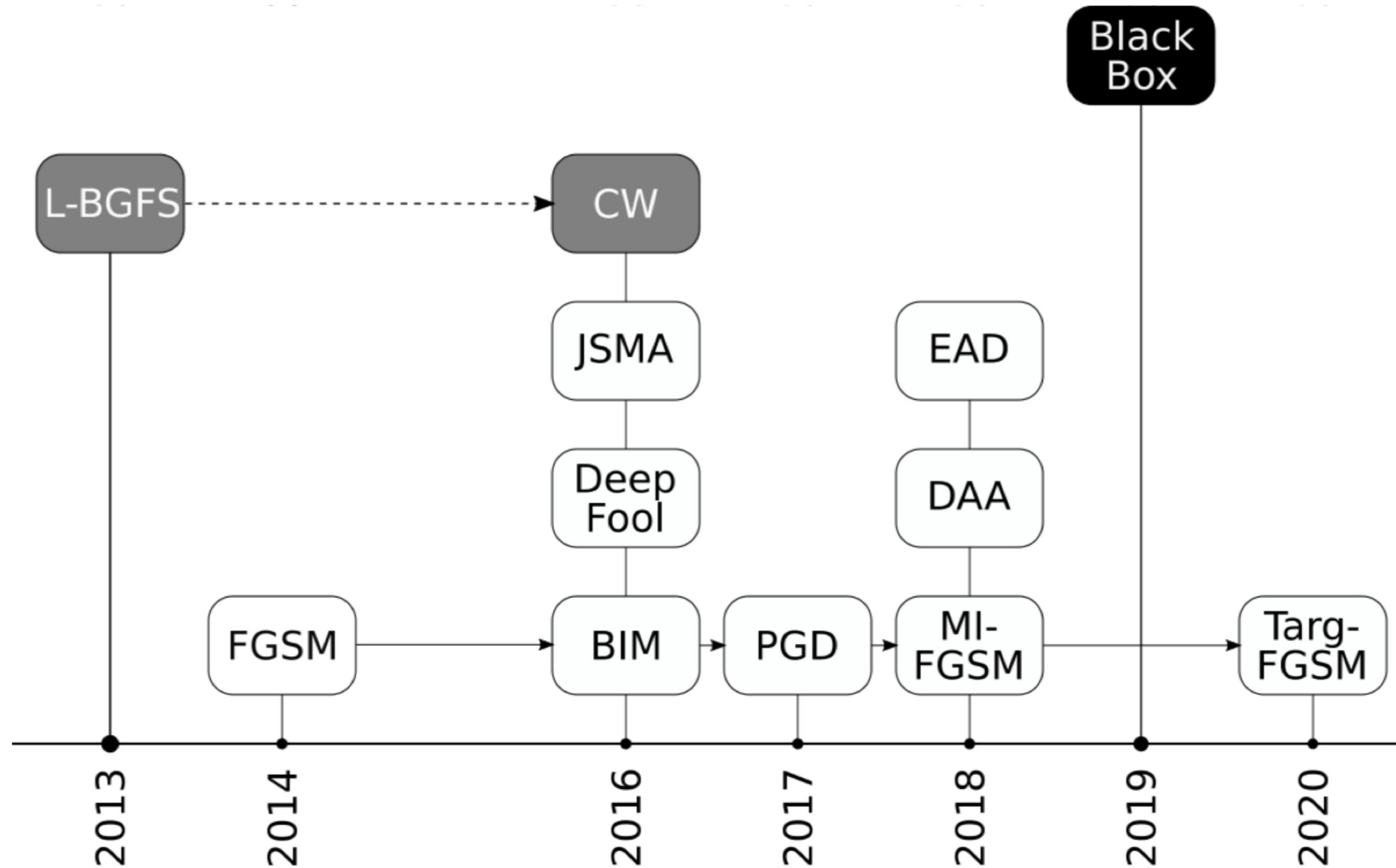


Architecture, loss

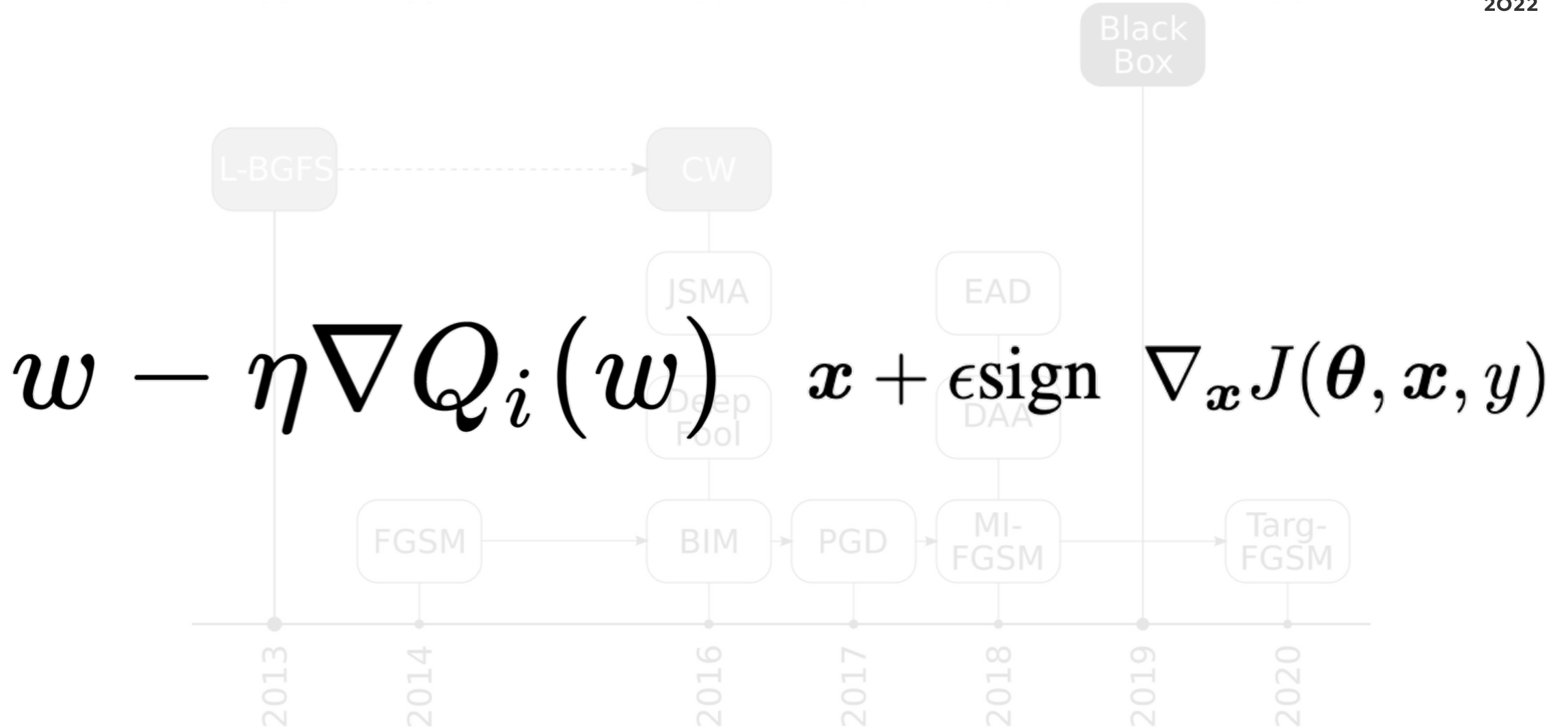


Everything

Adversarial examples generating methods



Adversarial examples generating philosophy



Frameworks

List selected for testing

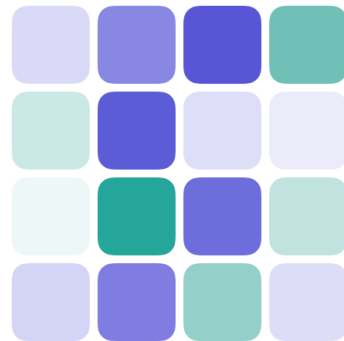


AdvBox

ART

Foolbox

DeepRobust



Adversarial
Robustness
Toolbox



Frameworks

AdvBox



AdvBox

ART

Foolbox

DeepRobust

- <https://github.com/advboxes/AdvBox>
- 2020
- Based on Foolbox v1
- It's alive!

Frameworks

AdvBox







AdvBox

ART

Foolbox

DeepRobust

BIM	FGSM	MI-FGSM	<u>DeepFool</u>
			

Frameworks

ART



Adversarial Robustness Toolbox

AdvBox

ART

Foolbox

DeepRobust

- <https://github.com/Trusted-AI/adversarial-robustness-toolbox>
- 2020
- Not bad

Frameworks

ART



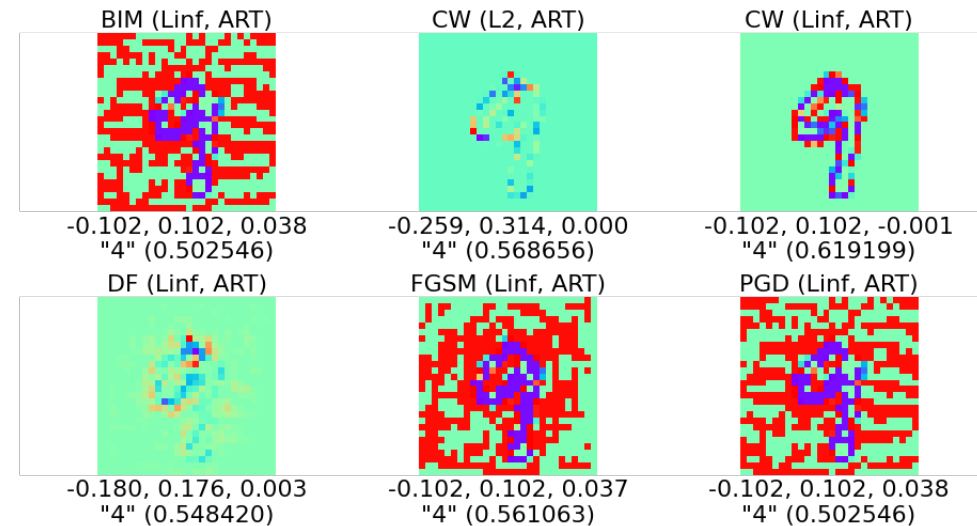
Adversarial Robustness Toolbox

AdvBox

ART

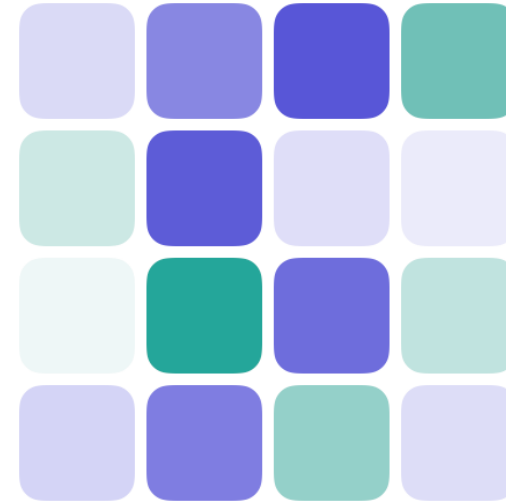
Foolbox

DeepRobust



Frameworks

Foolbox



AdvBox

ART

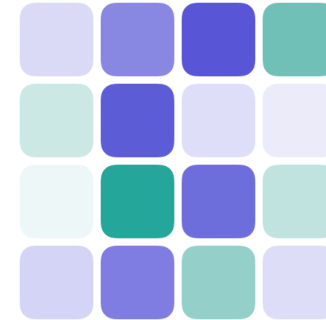
Foolbox

DeepRobust

- <https://github.com/bethgelab/foolbox>
- 2017
- Epic!

Frameworks

Foolbox

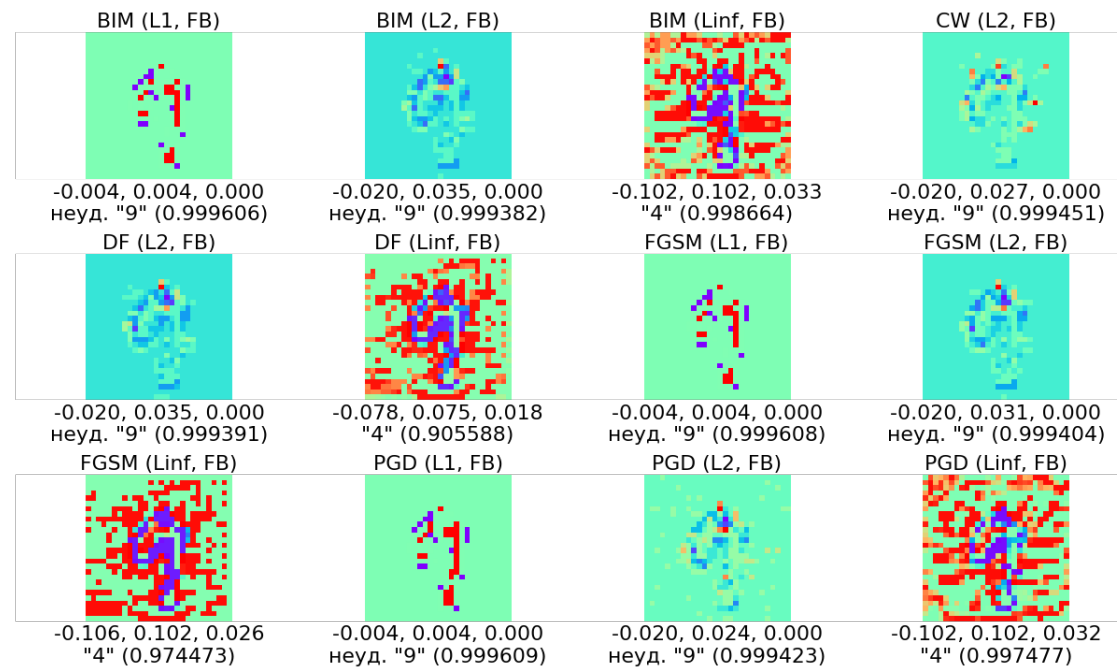


AdvBox

ART

Foolbox

DeepRobust



Frameworks

DeepRobust



AdvBox

ART

Foolbox

DeepRobust

- <https://github.com/DSE-MSU/DeepRobust>
- 2020
- not touching that one with a twenty-foot pole

Frameworks

DeepRobust



AdvBox

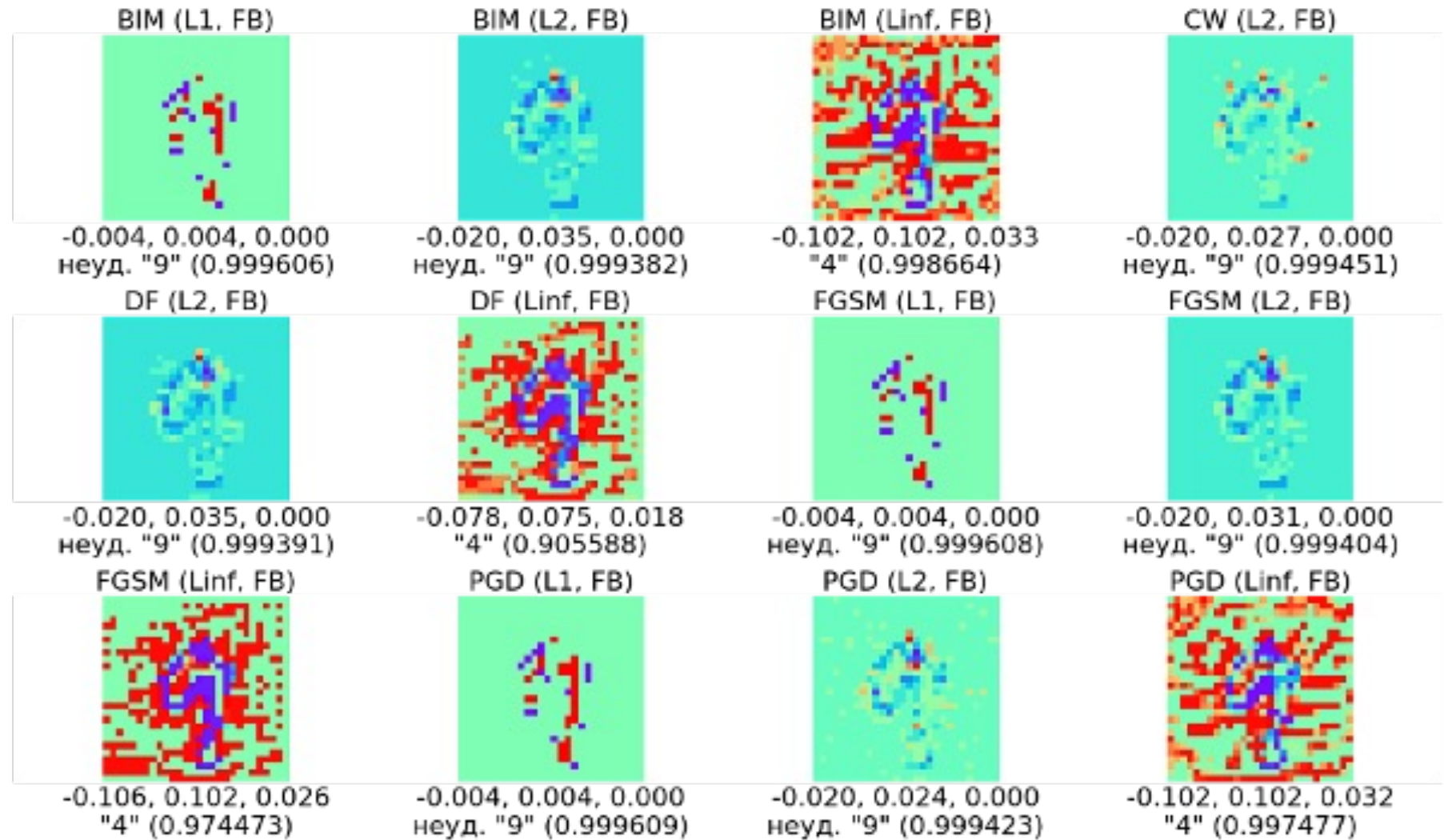
ART

Foolbox

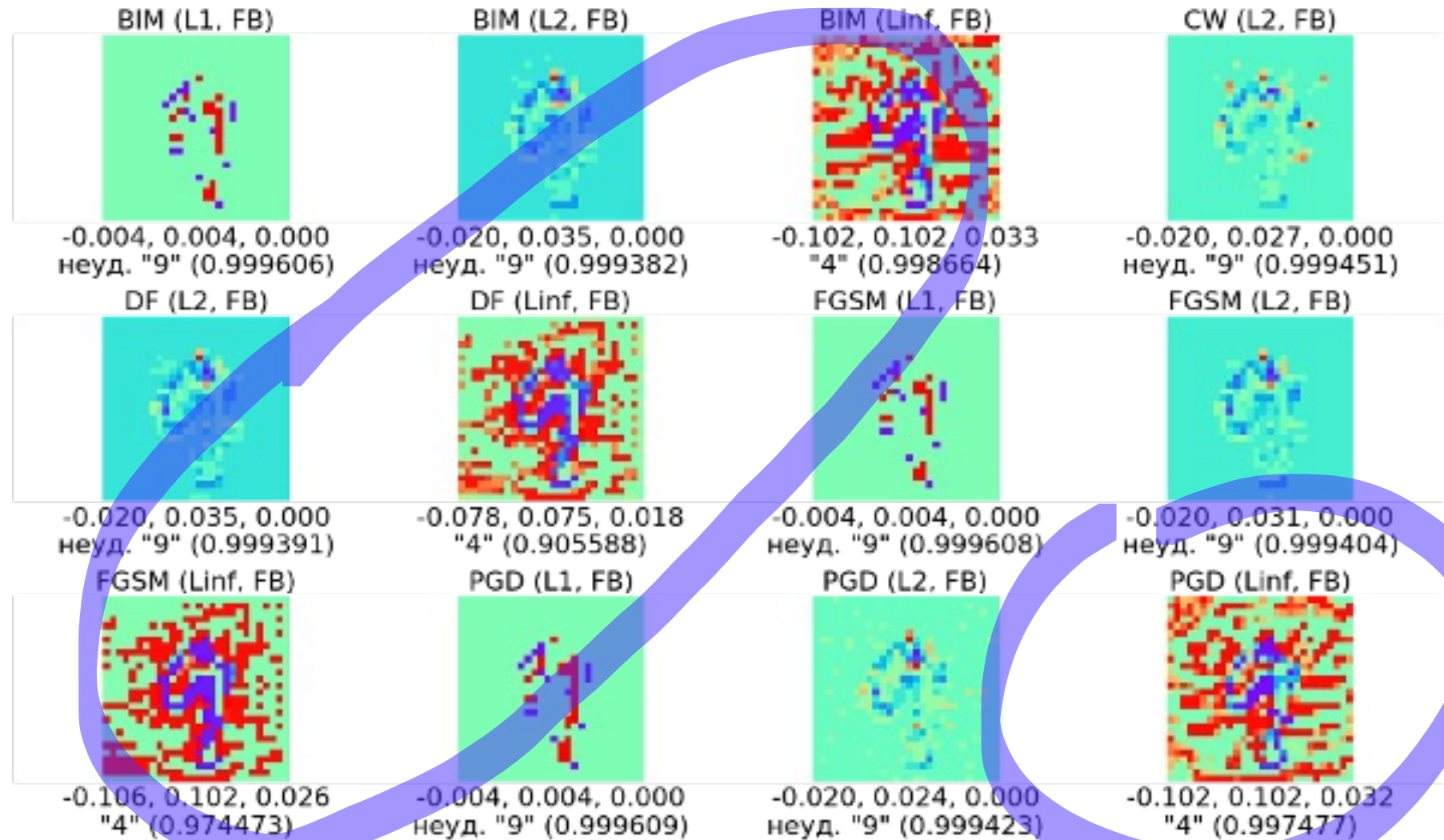
DeepRobust



Adversarial examples for class "9" obtained by Foolbox

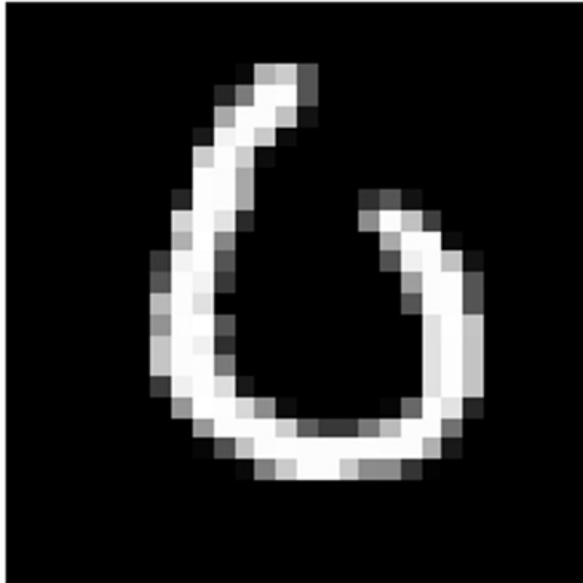


Adversarial examples for class "9" obtained by Foolbox



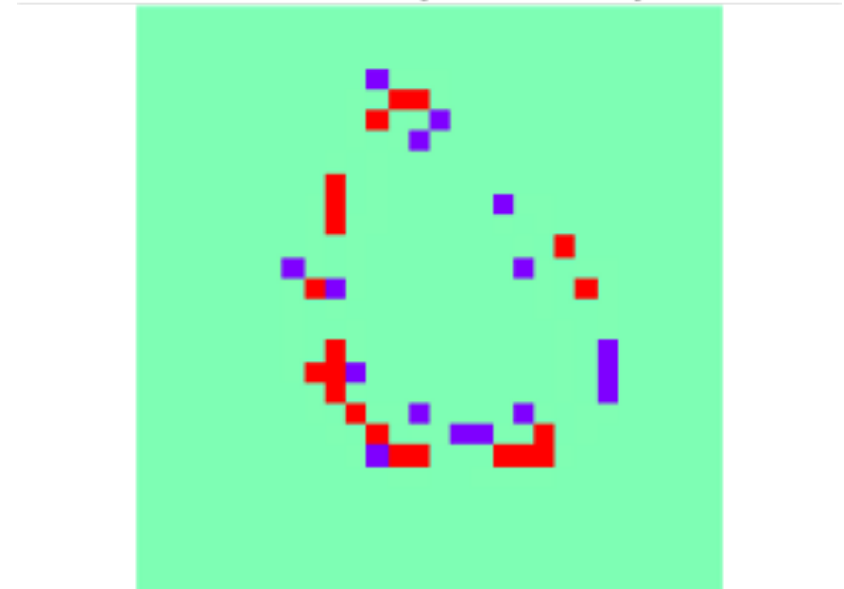
Perturbation value

“Adversarial example” for human



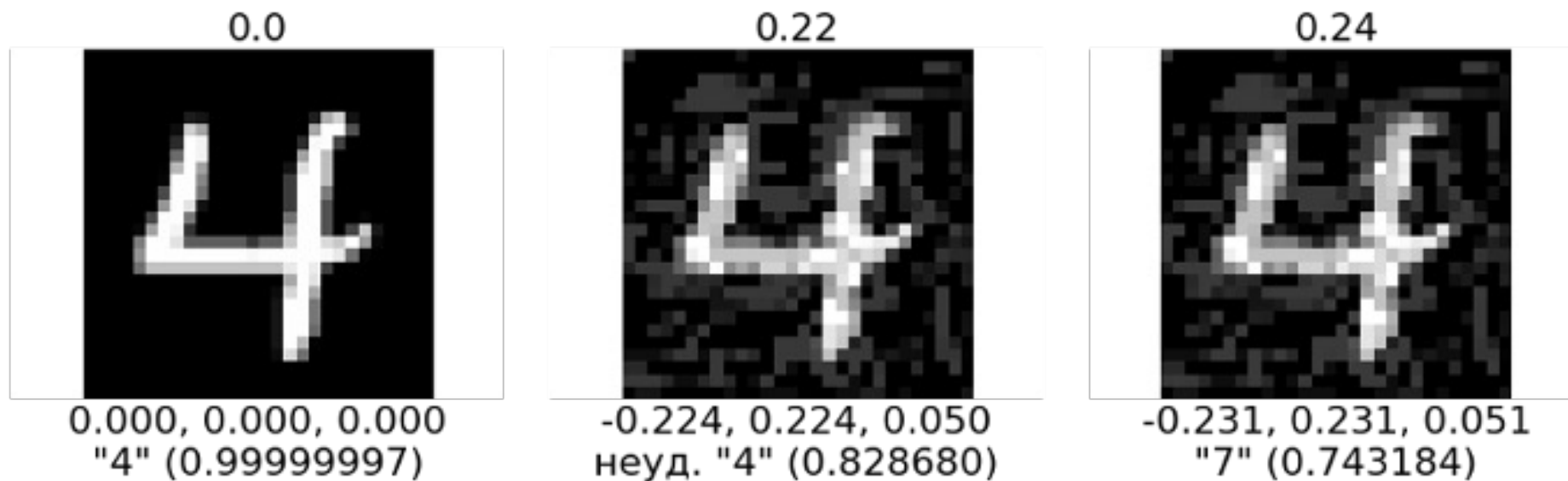
“0” 0.895

FGSM (L1, FB)



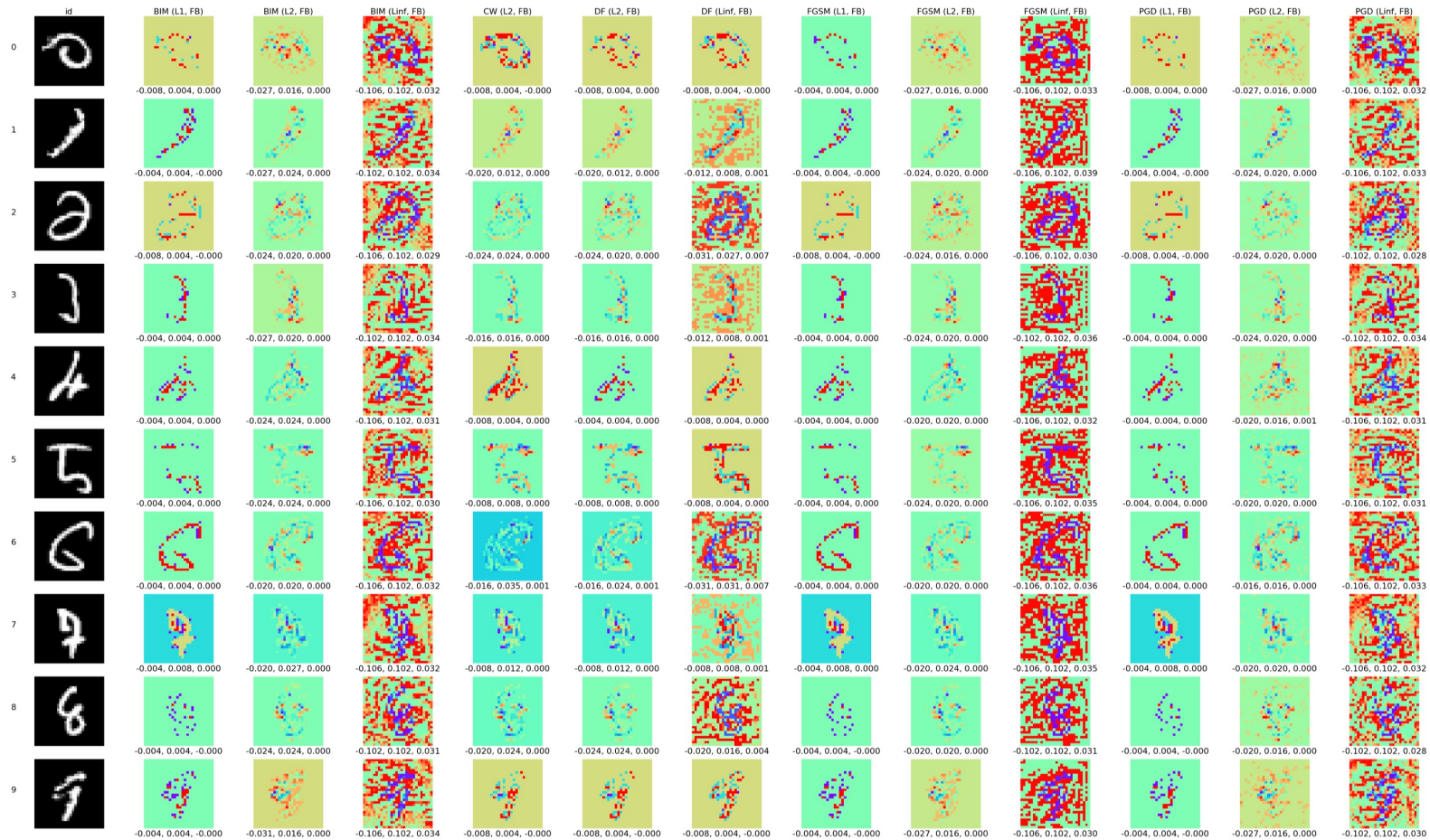
-0.004, 0.004, 0.000
“6” (0.893986)

Perturbation value



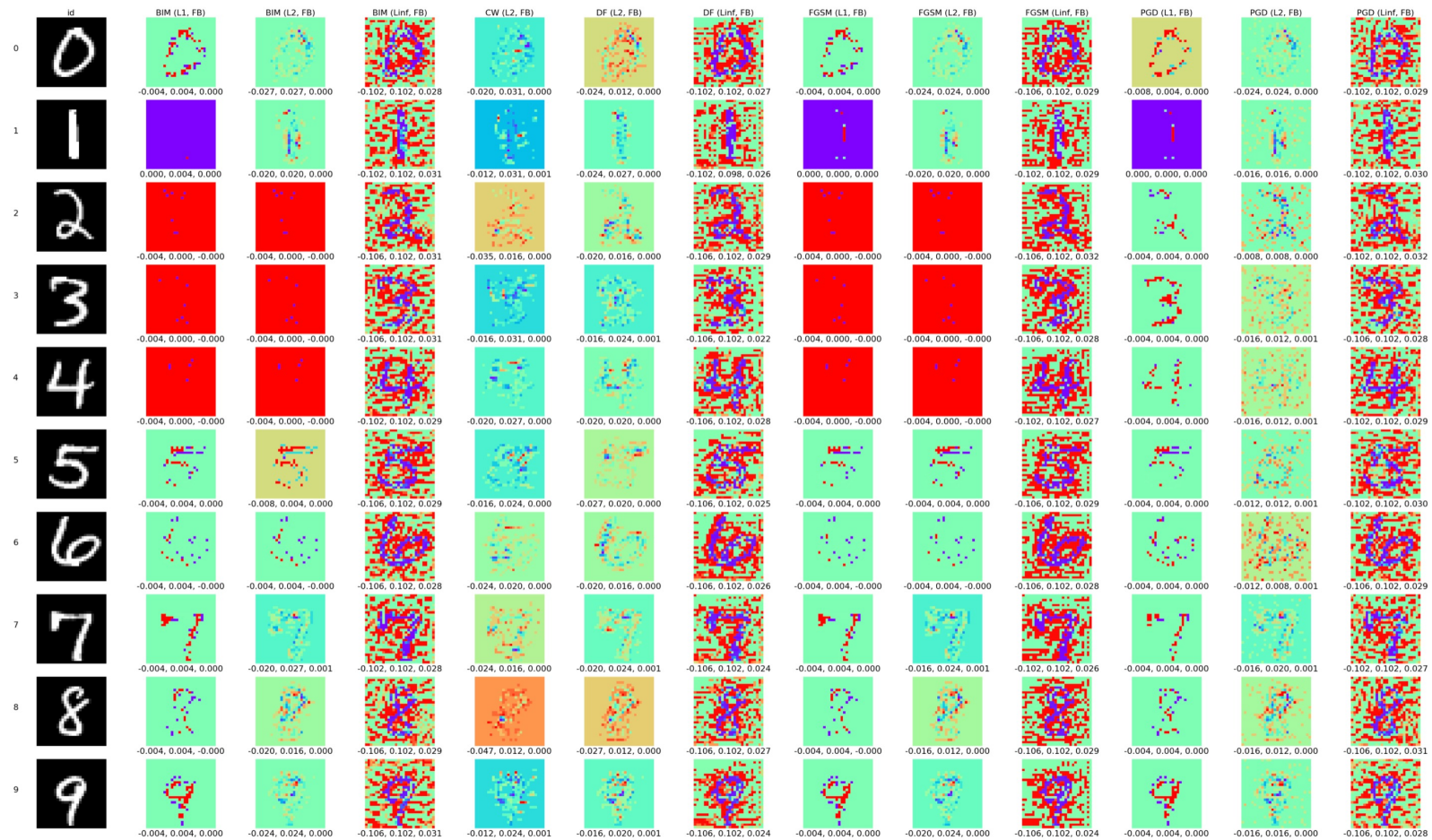
Perturbation value

Bad examples



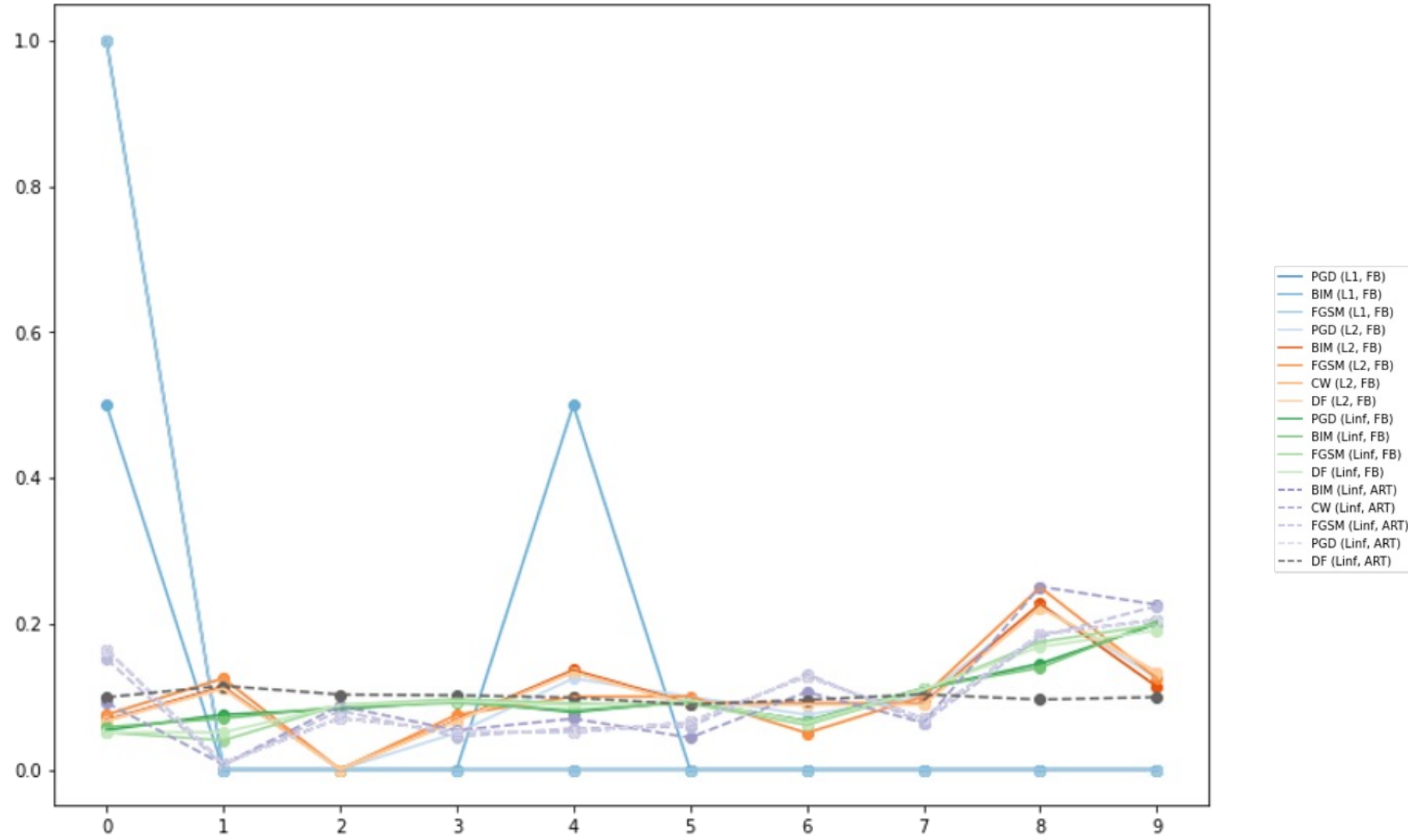
Perturbation value

Good examples



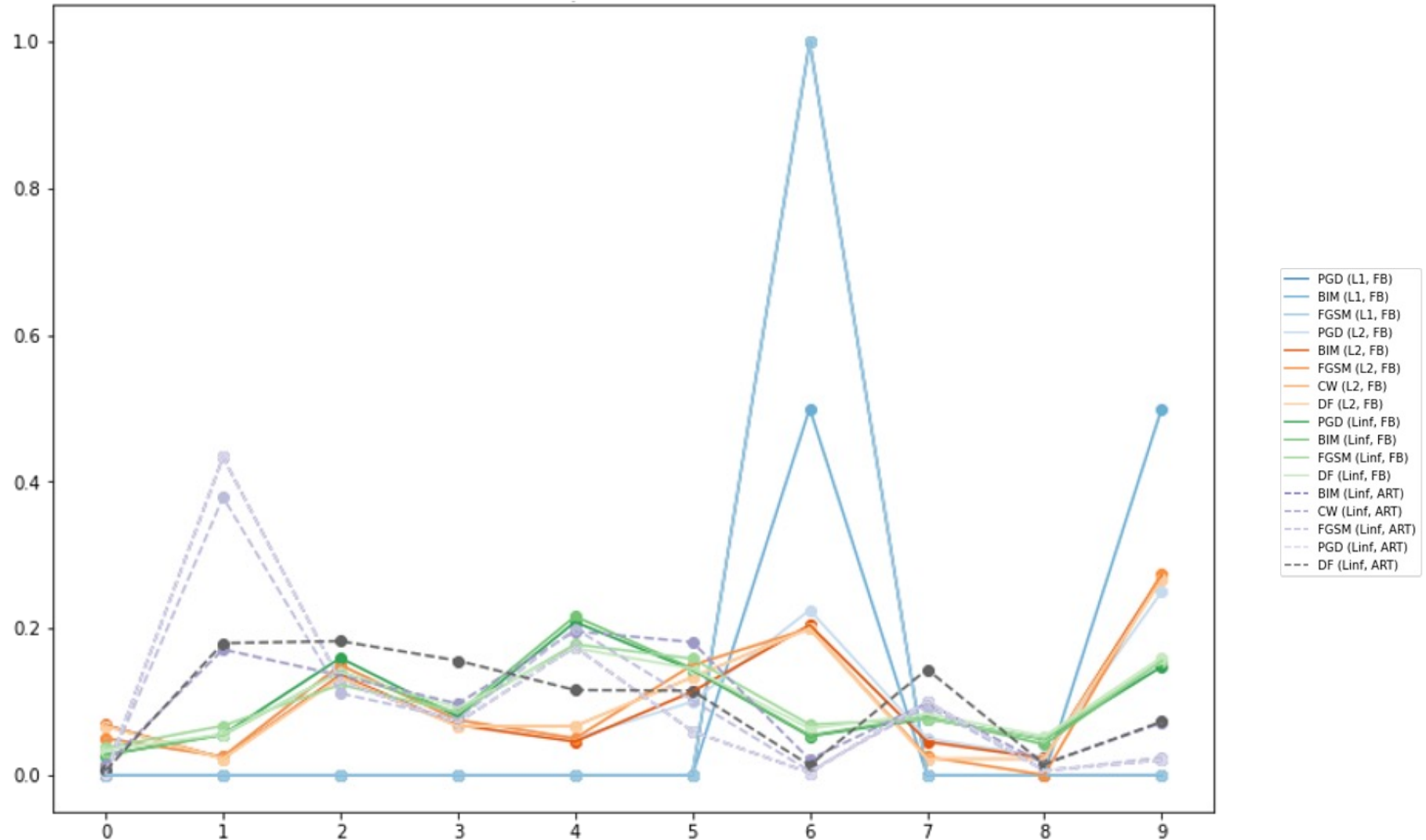
Attack performance vs label

Adversarial examples sources chart

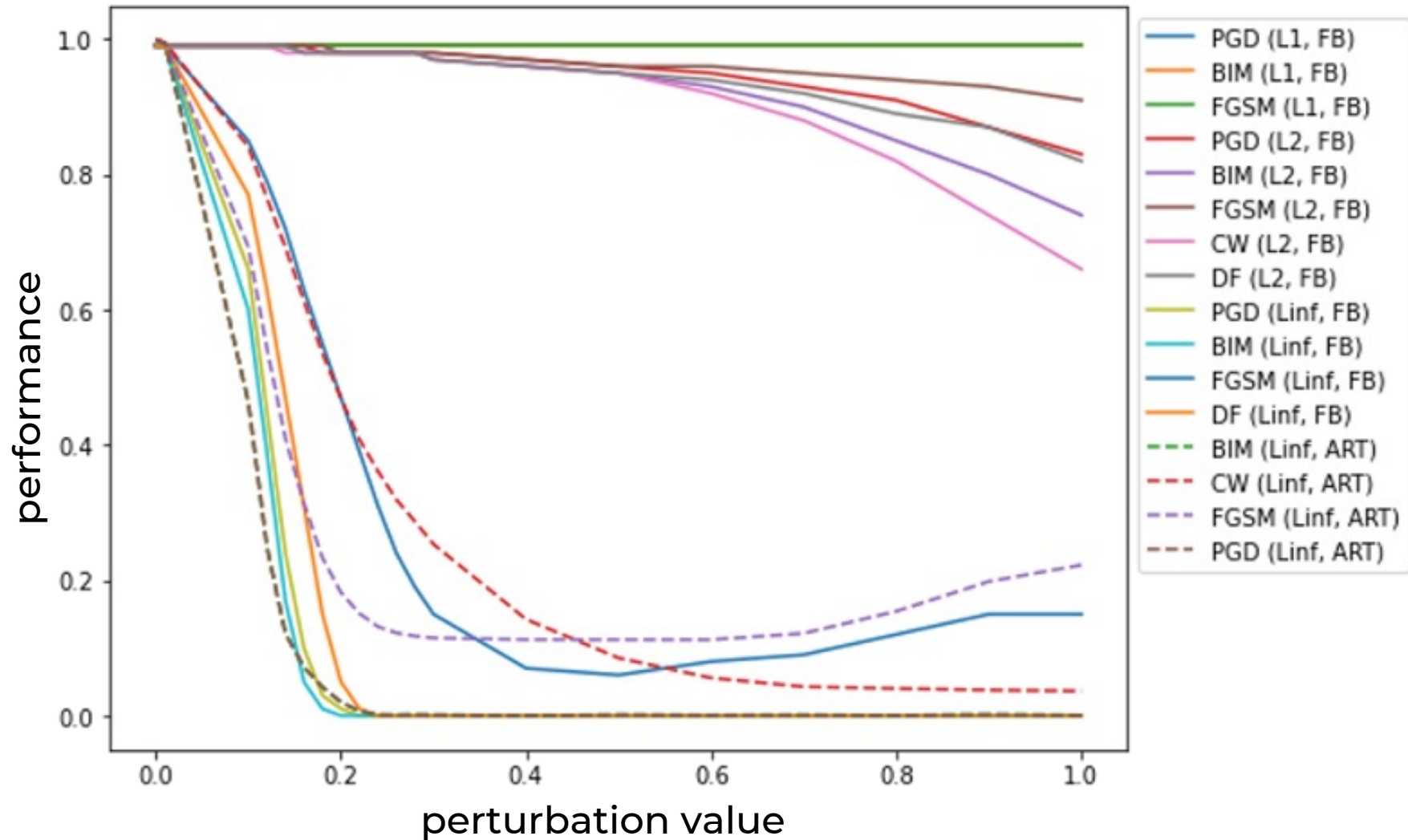


Attack performance vs label

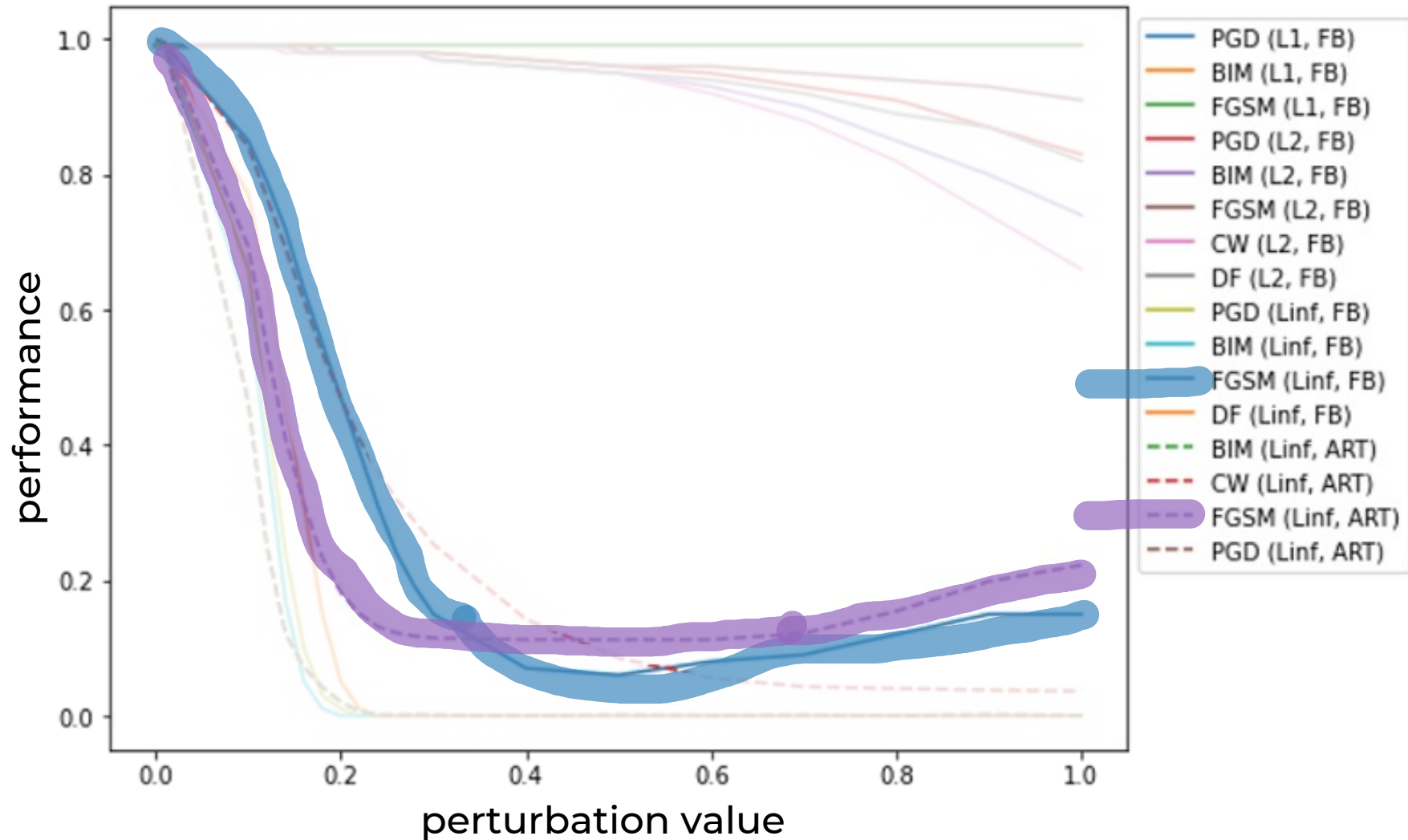
Adversarial examples targets chart



Benchmark of different methods implemented in ART and Foolbox



Benchmark of different methods implemented in ART and Foolbox



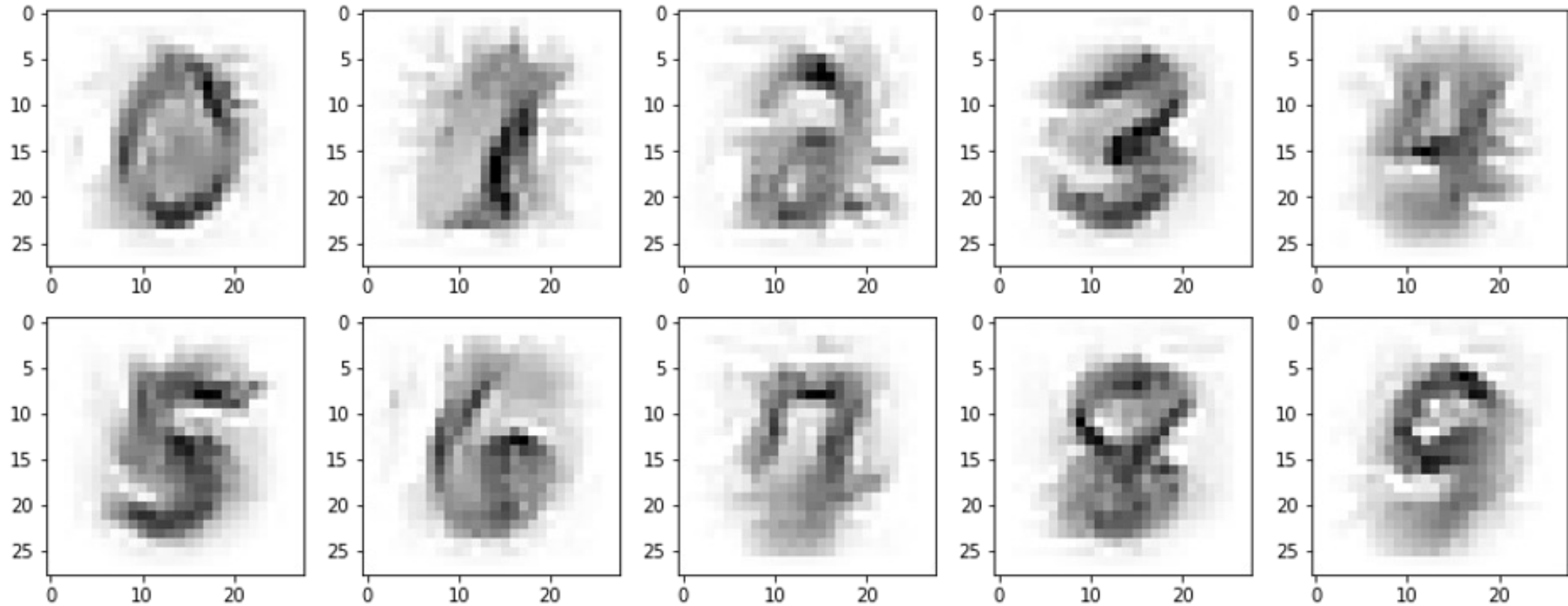
Beyond the scope

Black-box attacks



Beyond the scope

Data extraction



Thanks!



Danila Yu. Emelyanov
d.emelyanov@kryptonite.ru