



MAX-PLANCK-INSTITUT
FÜR BIOLOGISCHE KYBERNETIK



Pluribus One
seeing one in many



Pattern Recognition
and Applications Lab



University of
Cagliari, Italy

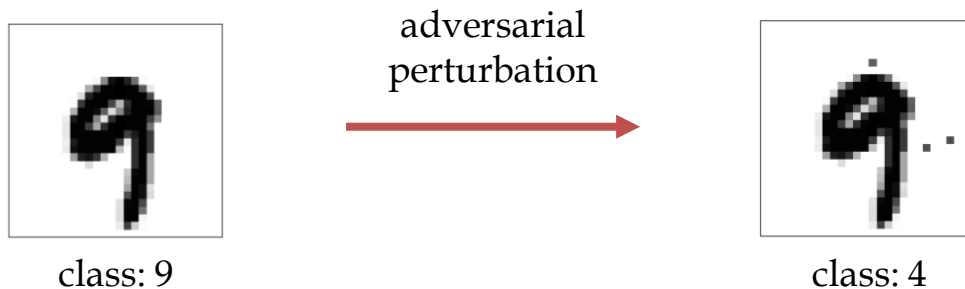
Fast minimum-norm adversarial attacks through adaptive norm constraints

Maura Pintor, Fabio Roli, Wieland Brendel, Battista Biggio

Thirty-fifth Conference on Neural Information Processing Systems

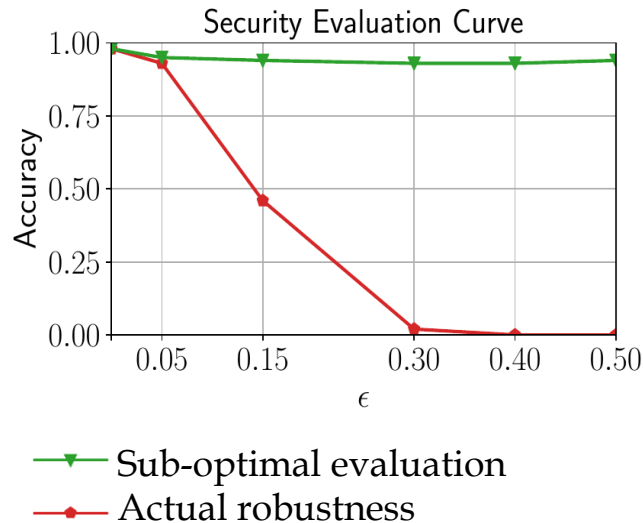


Introduction



$$\begin{aligned} \delta^* \in \arg \min_{\delta} \quad & \|\delta\|_p, \\ \text{s.t.} \quad & L(x + \delta, y, \theta) < 0, \\ & x + \delta \in [0, 1]^d, \end{aligned}$$

Challenges of evaluating adversarial robustness










The optimization depends on the **points** and on the **model** under attack




The risk of using sub-optimal hyperparameters might lead to **over-optimistic evaluations** (failing attacks)

Available algorithms do not usually maintain stable performances for different points

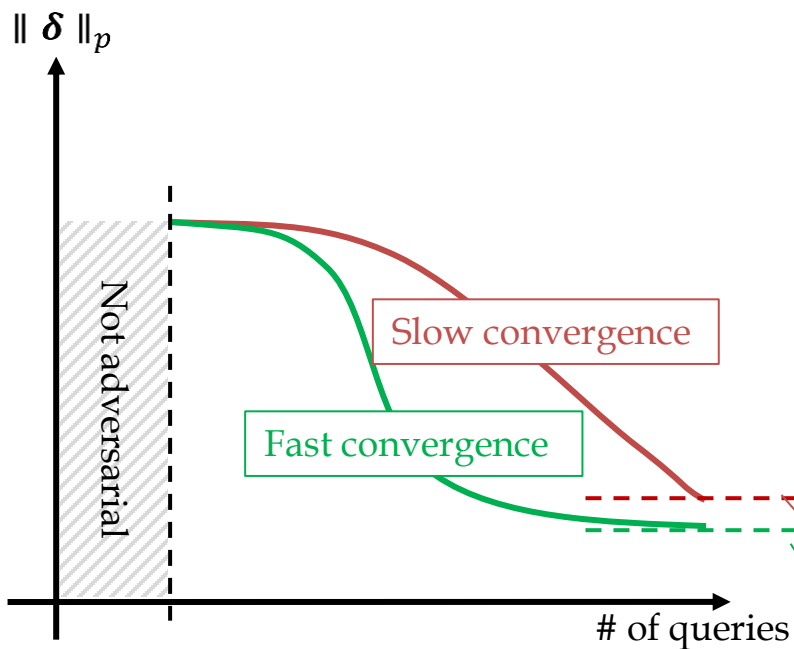
Most advanced attacks aim to obtain better solutions, at the cost of **longer execution**

Adversarial perturbations with minimum norm

- **Carlini-Wagner attack (CW)**
 - Requires many steps to converge 
- **Brendel&Bethge attack (BB)**
 - Needs initialization 
 - Suffers from poor initialization 
 - Complicated steps 
- **Fast Adaptive Boundary (FAB)**
 - Complicated steps 
 - Only untargeted version 
- **Decoupling Direction & Norm (DDN)**
 - Specific to L2 norm 

-  Long runtime
-  Sensitive to hyperparameters
-  Limited threat model

Our contributions



FMN



Fast convergence to good local optima



Works in different norms ($\ell_0, \ell_1, \ell_2, \ell_\infty$)



Easy tuning / robust to hyperparameter choice



Framework for evaluating existing attacks

Evaluation of minimum-norm attacks

4 different evaluation criteria

Bigger norm

Smaller norm

Our attack: Fast Minimum-Norm (FMN)

Algorithm 1 Fast Minimum-norm (FMN) Attack

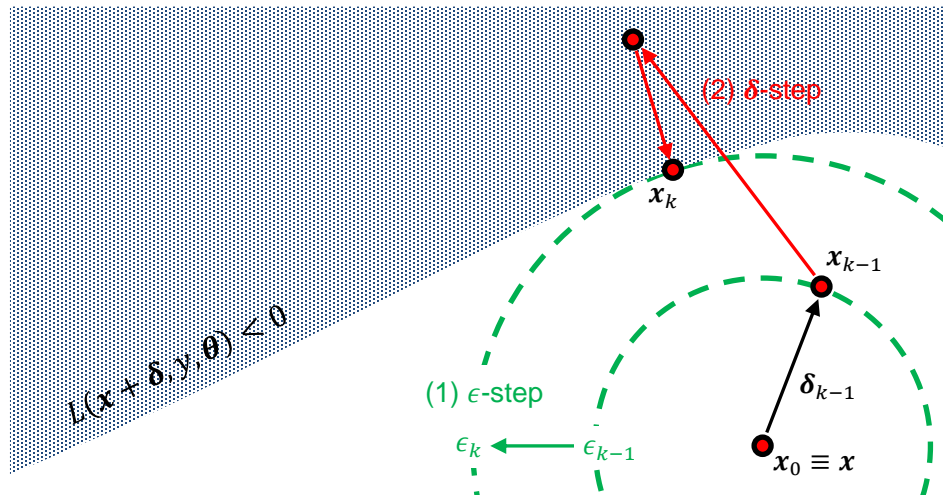
Input: \mathbf{x} , the input sample; t , a variable denoting whether the attack is targeted ($t = +1$) or untargeted ($t = -1$); y , the target (true) class label if the attack is targeted (untargeted); γ_0 and γ_K , the initial and final ϵ -step sizes; α_0 and α_K , the initial and final δ -step sizes; K , the total number of iterations.

Output: The minimum-norm adversarial example \mathbf{x}^* .

```

1:  $\mathbf{x}_0 \leftarrow \mathbf{x}, \epsilon_0 \equiv 0, \delta_0 \leftarrow \mathbf{0}, \delta^* \leftarrow \infty$ 
2: for  $k = 1, \dots, K$  do
3:    $\mathbf{g} \leftarrow t \cdot \nabla_{\delta} L(\mathbf{x}_{k-1} + \delta, y, \theta)$  // loss gradient
4:    $\gamma_k \leftarrow h(\gamma_0, \gamma_K, k, K)$  //  $\epsilon$ -step size decay (Eq. 7)
5:   if  $L(\mathbf{x}_{k-1}, y, \theta) \geq 0$  then
6:      $\epsilon_k = \|\delta_{k-1}\|_p + L(\mathbf{x}_{k-1}, y, \theta) / \|\mathbf{g}\|_q$  if adversarial not found yet else  $\epsilon_k = \epsilon_{k-1}(1 + \gamma_k)$ 
7:   else
8:     if  $\|\delta_{k-1}\|_p \leq \|\delta^*\|_p$  then
9:        $\delta^* \leftarrow \delta_{k-1}$  // update best min-norm solution
10:    end if
11:     $\epsilon_k = \min(\epsilon_{k-1}(1 - \gamma_k), \|\delta^*\|_p)$ 
12:  end if
13:   $\alpha_k \leftarrow h(\alpha_0, \alpha_K, k, K)$  //  $\delta$ -step size decay (Eq. 7)
14:   $\delta_k \leftarrow \delta_{k-1} + \alpha_k \cdot \mathbf{g} / \|\mathbf{g}\|_2$ 
15:   $\delta_k \leftarrow \Pi_{\epsilon}(\mathbf{x}_0 + \delta_k) - \mathbf{x}_0$ 
16:   $\delta_k \leftarrow \text{clip}(\mathbf{x}_0 + \delta_k) - \mathbf{x}_0$ 
17:   $\mathbf{x}_k \leftarrow \mathbf{x}_0 + \delta_k$ 
18: end for
19: return  $\mathbf{x}^* \leftarrow \mathbf{x}_0 + \delta^*$ 

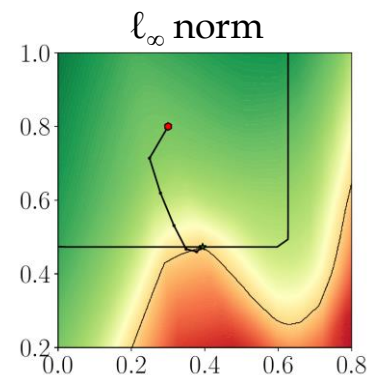
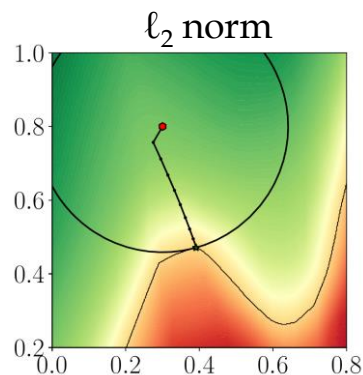
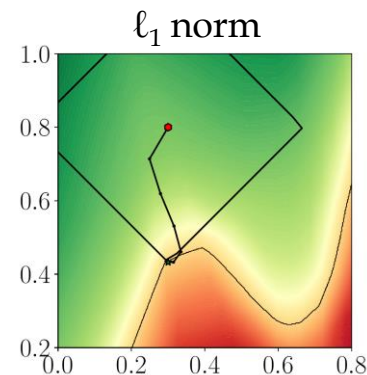
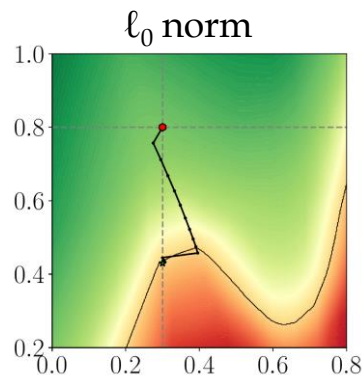
```



Perturbation models

FMN can find minimum-norm perturbations in 4 different ℓ_p norms

In each iteration, the attack performs a step in the direction of the gradient, and then projects the point back into the ℓ_p -ball of the corresponding norm

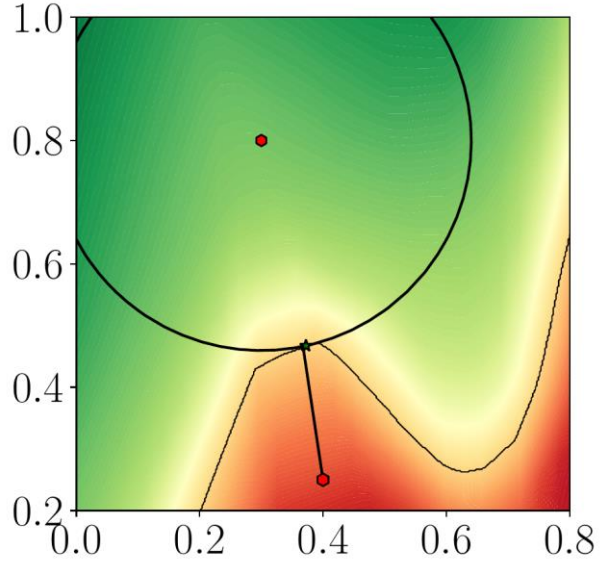


Adversarial Initialization

FMN can be initialized from a point in the target class

Finds the boundary quickly with an initial line search

Refines the results with the remaining iterations



Experimental setup

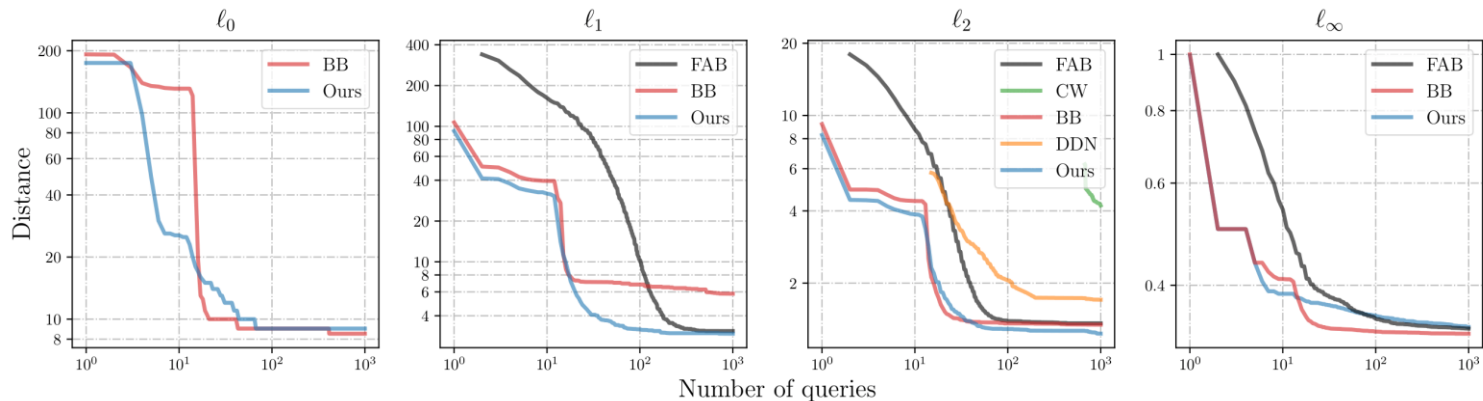
- 1 standard training model + 7 robust models in MNIST and CIFAR10 datasets (+2 ImageNet models)
- Comparison against 4 state-of-the-art minimum-distance attacks
- Targeted and untargeted scenario
- Evaluation across 4 ℓ_p distances ($p \in \{0, 1, 2, \infty\}$)

Evaluation

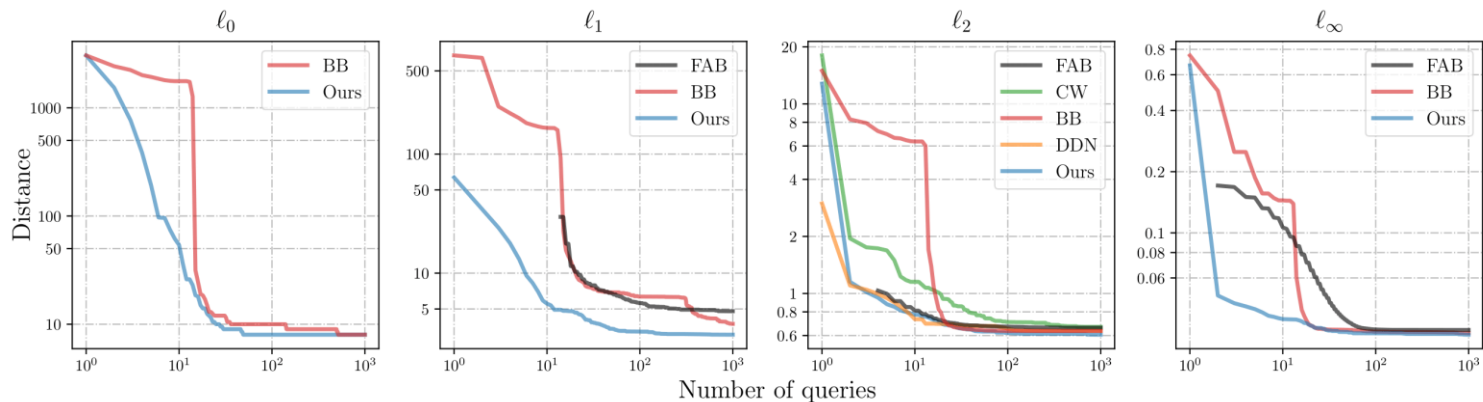
- Norm of perturbation
- Time-efficiency
- Query-efficiency
- Robustness to hyperparameters choice

Query-distortion curves

MNIST challenge



CIFAR challenge



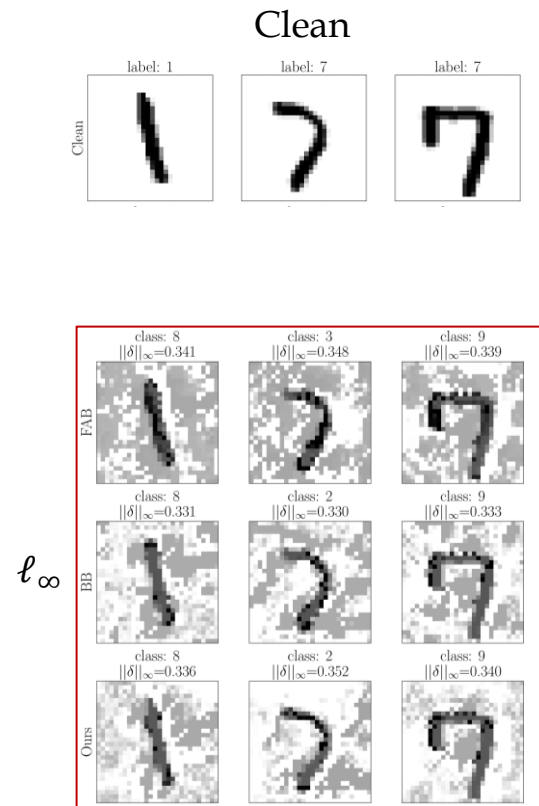
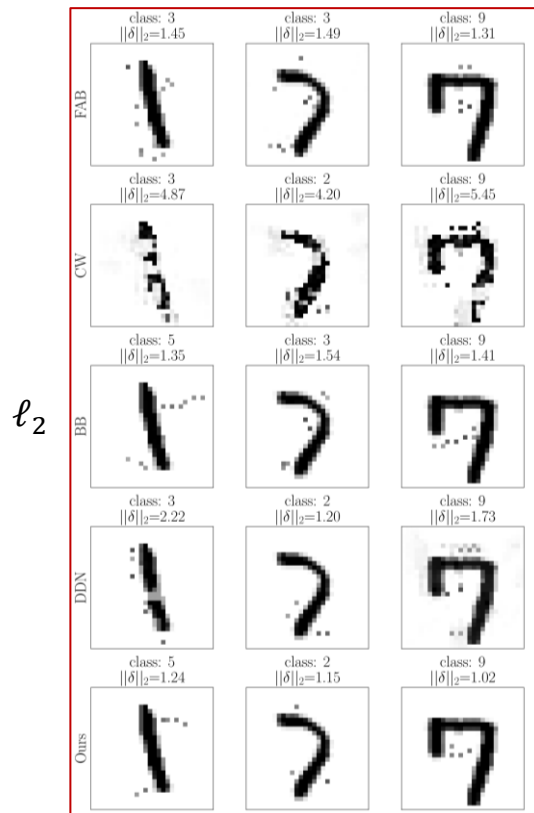
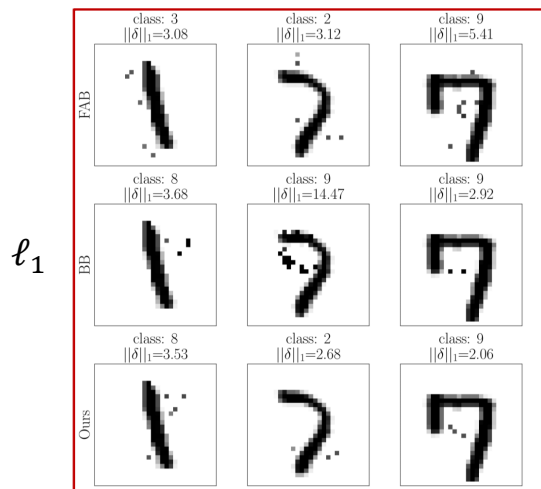
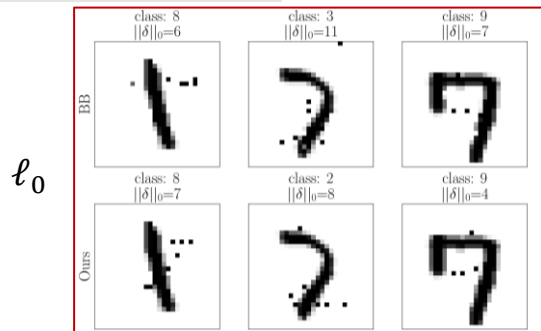
Time efficiency

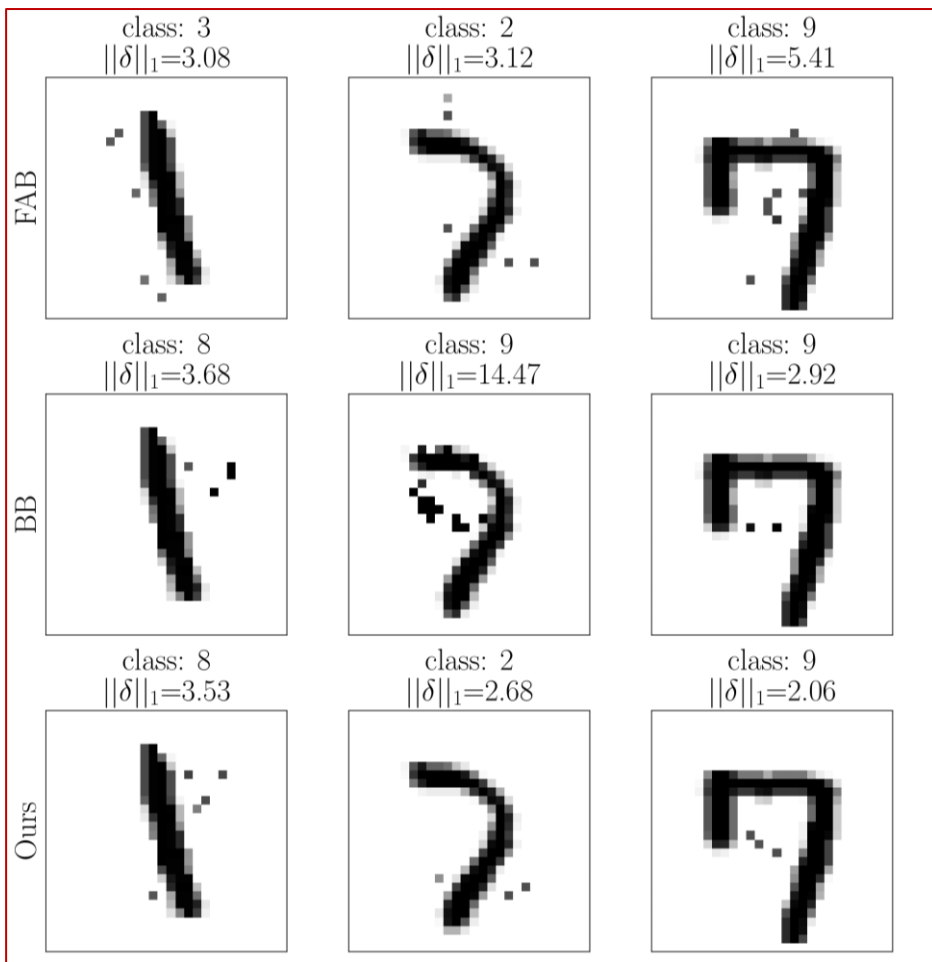
		Mnist Targeted	Mnist Untargeted	Cifar Targeted	Cifar Untargeted
l_0	BB	61.79 ± 0.79	10.93 ± 0.84	102.84 ± 2.85	49.08 ± 2.23
	Ours	5.24 ± 0.56	5.30 ± 0.51	29.07 ± 2.56	29.23 ± 2.58
l_1	FAB	—	10.29 ± 2.02	—	100.53 ± 14.28
	BB	43.49 ± 0.23	7.02 ± 0.29	72.12 ± 2.73	35.85 ± 2.85
	Ours	5.54 ± 0.50	5.56 ± 0.49	29.34 ± 2.89	29.90 ± 2.22
l_2	FAB	—	11.27 ± 1.89	—	100.86 ± 14.37
	CW	4.50 ± 0.56	4.50 ± 0.59	29.46 ± 3.19	29.51 ± 3.13
	BB	26.73 ± 0.52	4.54 ± 0.45	52.39 ± 3.16	30.12 ± 3.01
	DDN	3.69 ± 0.54	3.69 ± 0.54	27.58 ± 3.44	27.75 ± 3.13
	Ours	4.80 ± 0.58	4.79 ± 0.61	28.55 ± 2.74	28.39 ± 3.04
l_∞	FAB	—	11.83 ± 1.92	—	101.34 ± 14.48
	BB	37.96 ± 1.36	14.71 ± 1.48	86.75 ± 2.71	62.11 ± 0.68
	Ours	4.62 ± 0.60	4.62 ± 0.58	28.33 ± 3.08	28.33 ± 3.03

 Faster queries

 Easy tuning of the hyperparameters (more results in the paper)

Adversarial examples





Conclusions

Fast-Minimum Norm Attack

- Works in different norms
- Comparable or better norm of perturbation
- Query- and time-efficient
- Robust to hyperparameter choices

We provide

- Extensive experiments
- Open-source code

Available implementations:

- <https://github.com/pralab/Fast-Minimum-Norm-FMN-Attack>
- <https://github.com/bethgelab/foolbox>
- <https://github.com/jeromerony/adversarial-library>
- <https://github.com/pralab/secml>





MAX-PLANCK-INSTITUT
FÜR BIOLOGISCHE KYBERNETIK



Pluribus One
seeing one in many



Pattern Recognition
and Applications Lab



University of
Cagliari, Italy

**Fast minimum-norm adversarial attacks through
adaptive norm constraints**

Thank you for listening!