

# Adversarial attacks via Sequential Quadratic Programming

Théo Beuzeville<sup>1,3</sup> Pierre Boudier<sup>2</sup> Alfredo Buttari<sup>3,5</sup> Serge Gratton<sup>4</sup> Theo  
Mary<sup>5,6</sup> Stéphane Pralet<sup>1</sup> Erkan Ulker<sup>1</sup>

<sup>1</sup>ATOS

<sup>2</sup>NVIDIA

<sup>3</sup>IRIT

<sup>4</sup>Toulouse INP-ENSEEIH

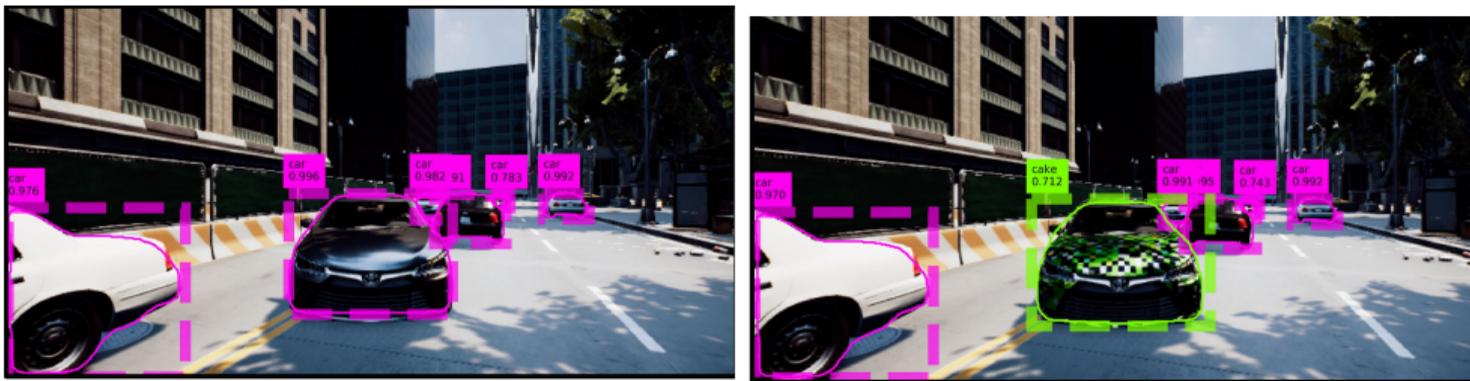
<sup>5</sup>CNRS

<sup>6</sup>LIP6

Sparse Days, June 2022

# Adversarial attacks

The approach consists in computing the **smallest norm perturbation** on input data such that, for a given input  $x$ , the perturbed input is **misclassified** by the neural network, that is, it erroneously affects the input to a given class  $j$  instead of the expected one.



**Figure:** The camouflage fools the Mask R-CNN object detector (on the right), whereas plain colors (on the left) is being correctly detected.

Mathematically, the adversarial perturbation is defined as the solution of the following minimization problem:

$$\begin{aligned} & \text{Solve} \\ & \min \|\Delta x\|^2 \\ & \text{subject to} \\ & C(x + \Delta x) = j. \end{aligned} \tag{1}$$

Where  $j$  is the target class and  $C(x + \Delta x)$  the class of the perturbed image.

This problem being difficult to solve, most of the approaches that generate adversarial examples commonly resort to solving the following problem:

**Solve**

$$\min \mathbf{c} \|\Delta \mathbf{x}\|^2 + \mathcal{L}(\mathbf{x} + \Delta \mathbf{x}, j), \quad (2)$$

where  $\mathcal{L}$  is a given loss of the image with respect to a given target class.



**Solve**

$$\min \|\Delta x\|^2$$

**subject to**

$$C(x + \Delta x) = j. \tag{3}$$

We introduce a novel method to try and solve this problem **more accurately** using **second order** information.

Thank you for your attention and see you at the poster corner.