



Adversarial Machine learning

Arman Akbari, Arash Akbari



x

“panda”

57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$

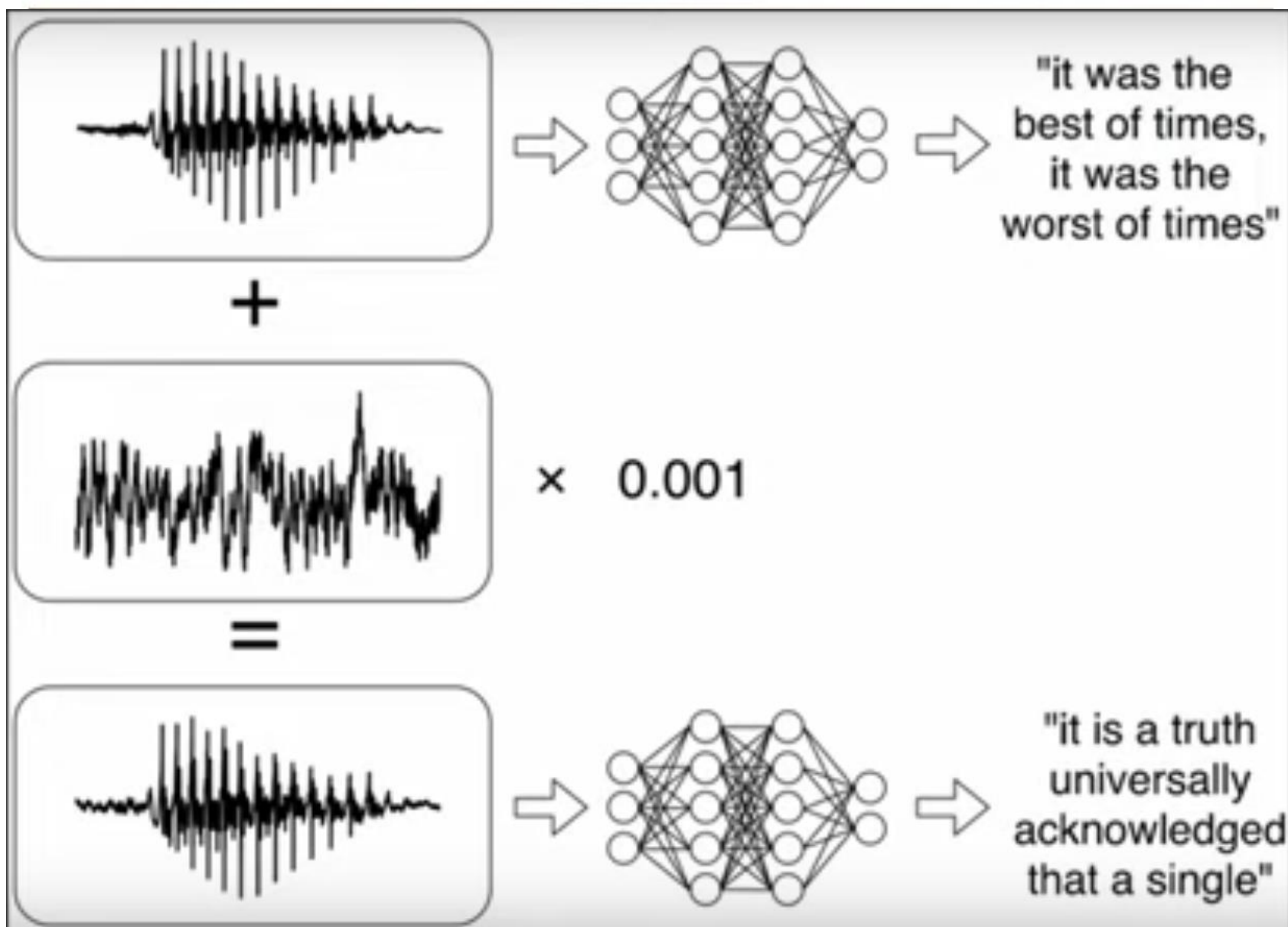
“nematode”

8.2% confidence

$=$



$x +$
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$
 “gibbon”
 99.3 % confidence



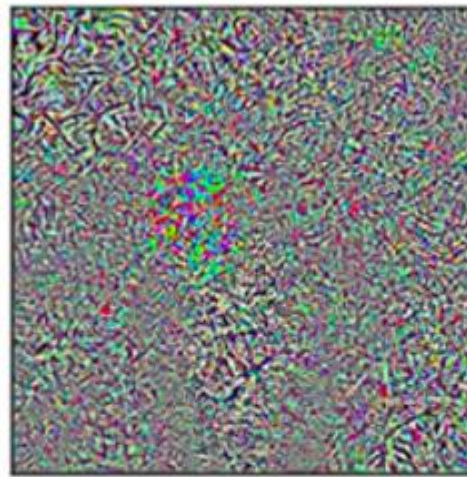
FGSM(Fast Gradient Sign Method)

- The linear view of adversarial examples suggests a fast way of generating them

$$\eta = \epsilon \text{sign}(\nabla_x J(\theta, x, y)).$$



+



=



X

97.3% macaw

$\text{sign}(\nabla_x J(\theta, X, Y))$

$X + \epsilon \cdot \text{sign}(\nabla_x J(\theta, X, Y))$

88.9% bookcase

Max-loss from



$$\begin{aligned} & \max_{\mathbf{x}'} \ell(\mathbf{y}, f_{\boldsymbol{\theta}}(\mathbf{x}')) \\ \text{s. t. } & d(\mathbf{x}, \mathbf{x}') \leq \varepsilon, \quad \mathbf{x}' \in [0, 1]^n, \end{aligned} \quad (1)$$

Robustness Radius



$$\begin{aligned} & \min_{\mathbf{x}'} d(\mathbf{x}, \mathbf{x}') \\ \text{s. t. } & \max_{i \neq y} f_{\boldsymbol{\theta}}^i(\mathbf{x}') \geq f_{\boldsymbol{\theta}}^y(\mathbf{x}'), \quad \mathbf{x}' \in [0, 1]^n, \end{aligned} \quad (2)$$

Adversarial Robustness



$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \max_{\mathbf{x}' \in \Delta(\mathbf{x})} \ell(\mathbf{y}, f_{\boldsymbol{\theta}}(\mathbf{x}')) \quad (3)$$