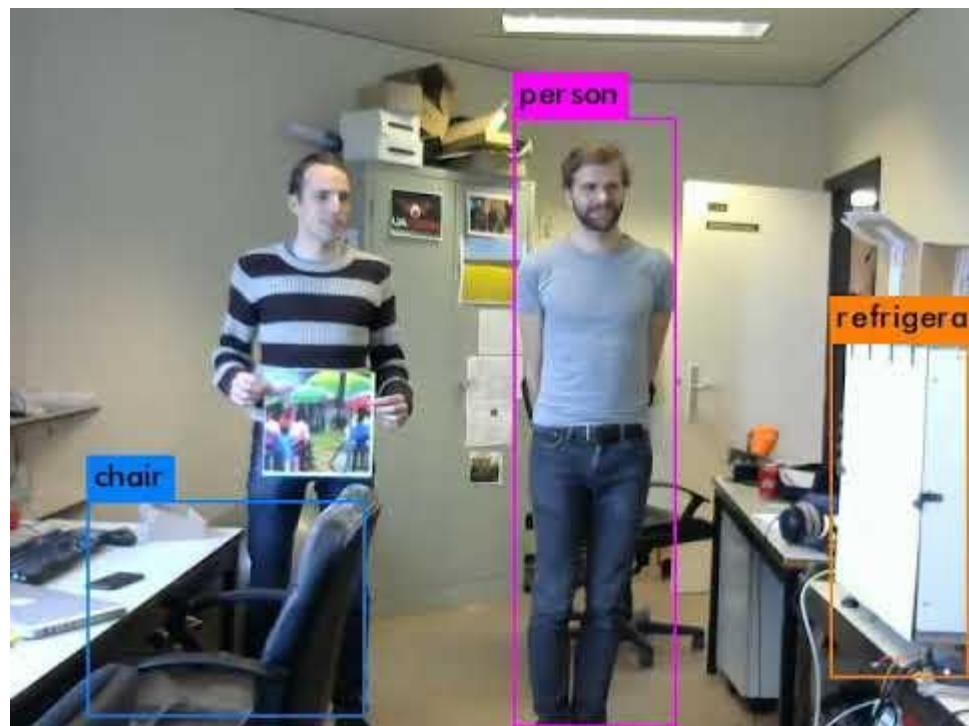


Adversarial attack in action







Audio

Transcription by Mozilla DeepSpeech



“without the dataset the article is useless”



“okay google browse to evil dot com”

https://nicholas.carlini.com/code/audio_adversarial_examples/

Tools

<https://github.com/ashishcse0031/Adversarial-Machine-Learning>

<https://github.com/tensorflow/cleverhans>

<https://github.com/IBM/adversarial-robustness-toolbox>

<https://github.com/bethgelab/foolbox>

Sign of the function tell's weather ,there is need to increase /decrease pixel value by a very small value ϵ to ensure that we do not go too far on the loss function surface and that the perturbation will be imperceptible.

$$J(\tilde{x}, \theta) \approx J(x, \theta) + (\tilde{x} - x)^\top \nabla_x J(x).$$

Maximize

$$J(x, \theta) + (\tilde{x} - x)^\top \nabla_x J(x)$$

subject to

$$\|\tilde{x} - x\|_\infty \leq \epsilon$$

$$\Rightarrow \tilde{x} = x + \epsilon \text{sign}(\nabla_x J(x)).$$



x

“panda”

57.7% confidence

+ .007 ×



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

“nematode”

8.2% confidence

=



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

99.3 % confidence