# Hacking Machine Learning Systems

**Cédric Simal**

# Motivation

**The current state of computer security**



http://gunshowcomic.com/648

# Plan

0. Why you should care
1. Security Framework
2. Common Attack types
3. Defenses

# All software can be hacked



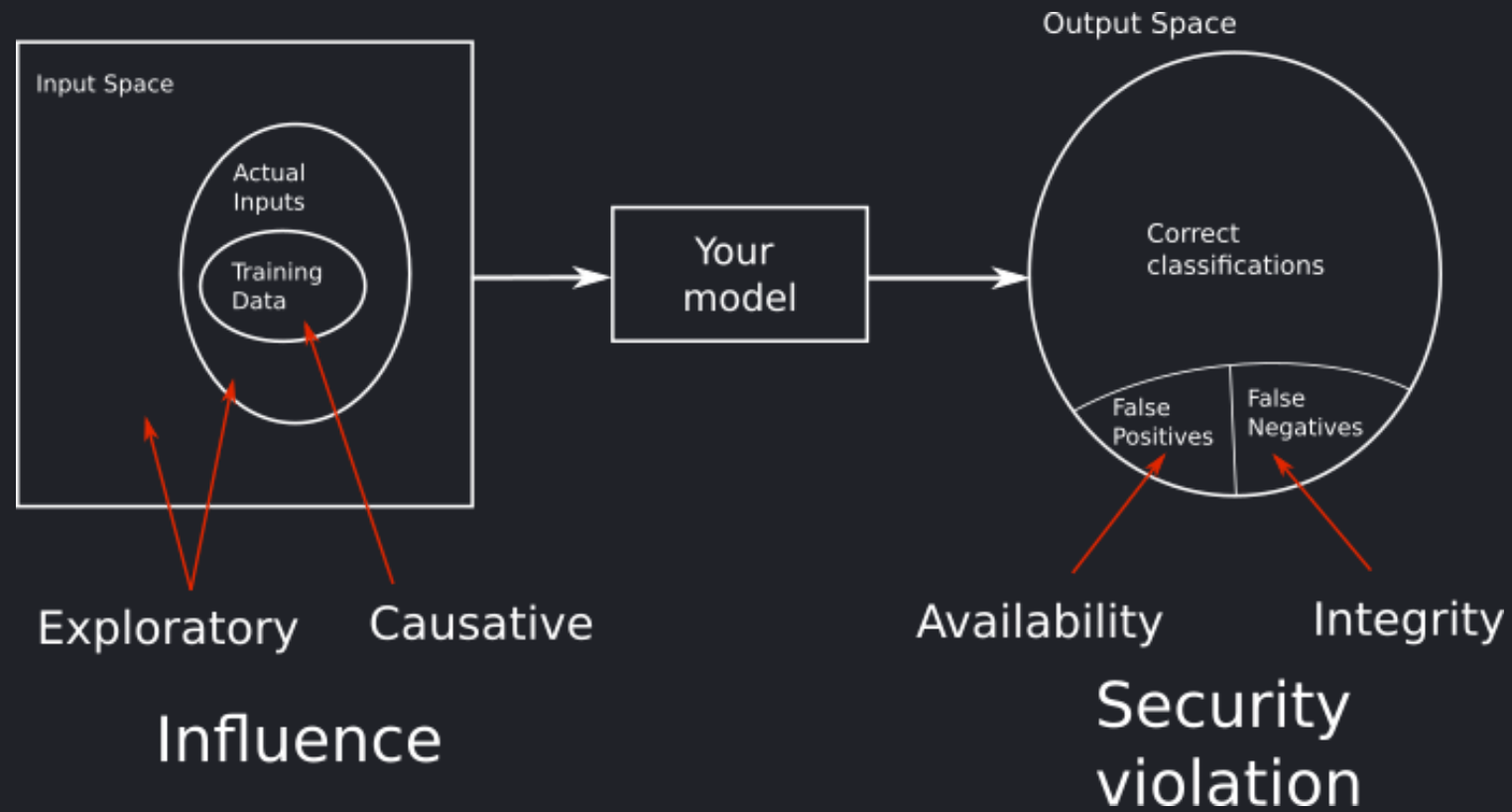https://knowyourmeme.com/memes/sites/tay-ai
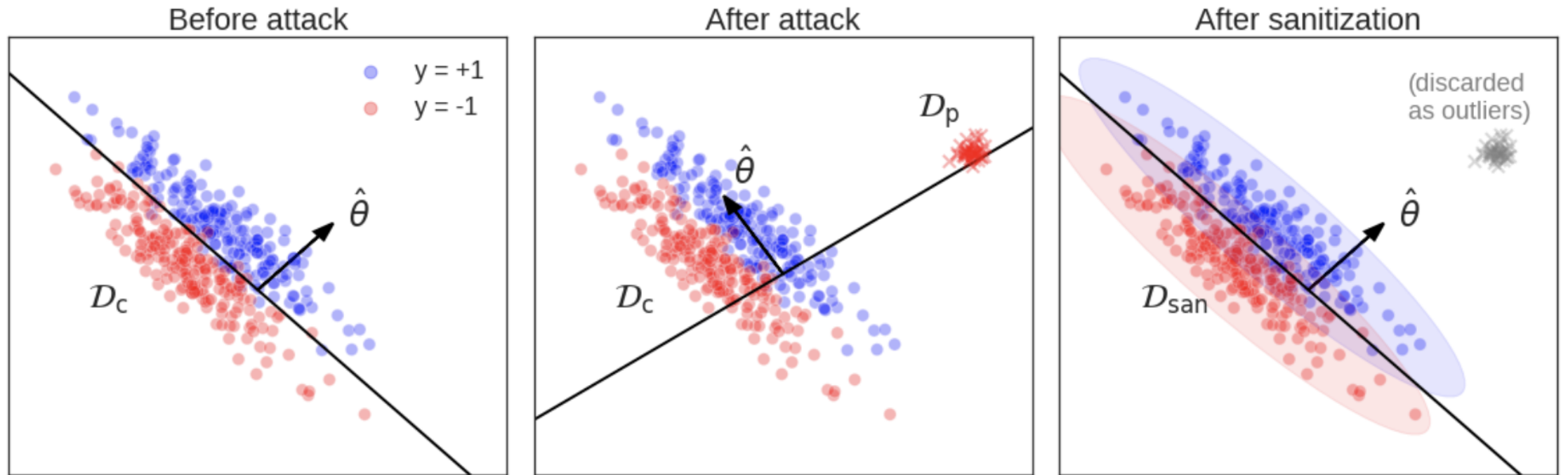
# Why attack ML?

- Degrade model performance
- Find misclassified inputs
- Steal model parameters
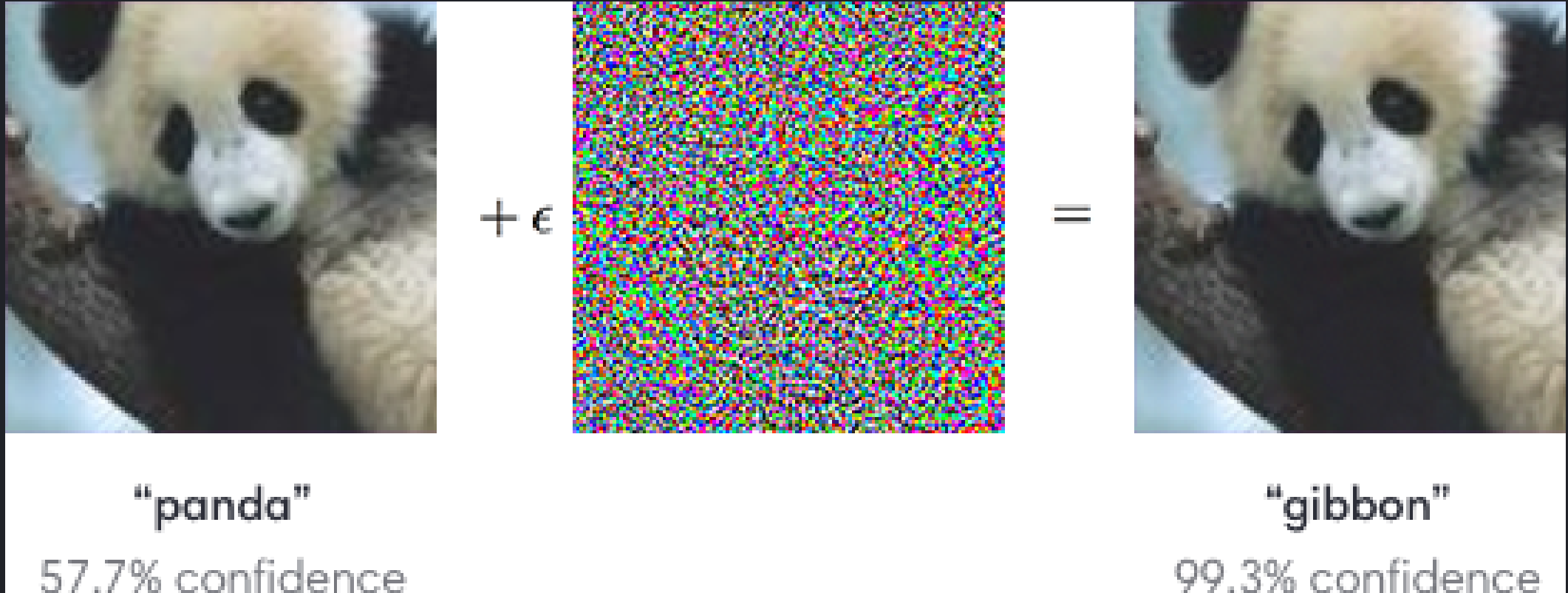- Steal training data

# Framework for attacks on ML



[Barreno et al.]

# Data Poisoning



[Koh et al.]

# Adversarial Examples



"panda"
57.7% confidence

$+\epsilon$
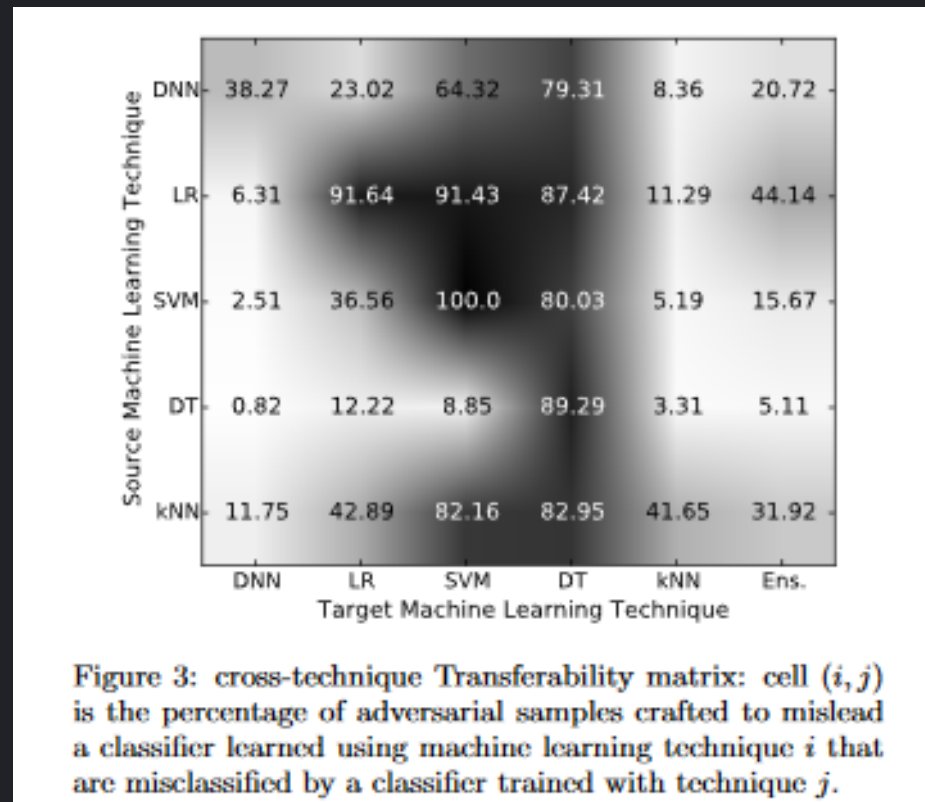
$=$
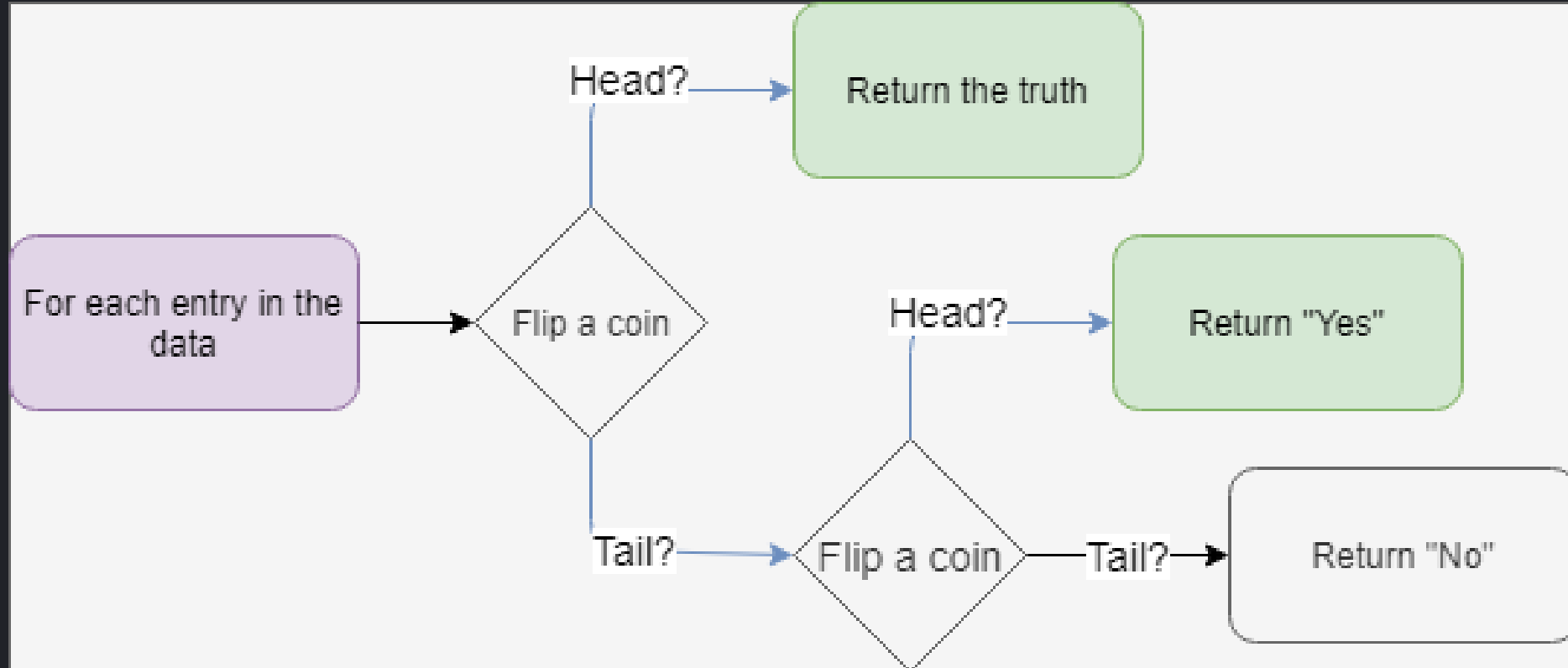
"gibbon"
99.3% confidence

openai.com

# Finding Adversarial Examples

- Whitebox : Fast Gradient Sign Method [Goodfellow et al.]
- Blackbox: Surrogate model [Papernot et al.]

# Attack Transferability



Figure 3: cross-technique Transferability matrix: cell $(i, j)$ is the percentage of adversarial samples crafted to mislead a classifier learned using machine learning technique $i$ that are misclassified by a classifier trained with technique $j$.

[Papernot et al.]

# Differential Privacy



towardsdatascience.com

# Defense Strategies

- Reject On Negative Impact (RONI)
- Adversarial Training
- Ensemble methods

# Take home message

- Review your data
- Check your inputs
- Include Security at the design stage

# Learn more!

Slides and sources at
https://github.com/csimal/ML-Hacking-Presentation