# Deep Deception
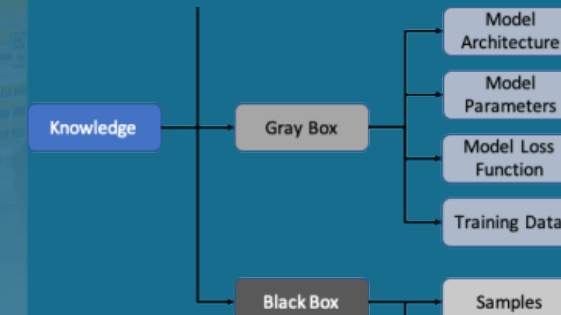
D. Farley, Z. Gastelum, T. Shead

**Artificial intelligence and machine learning for IAEA safeguards**

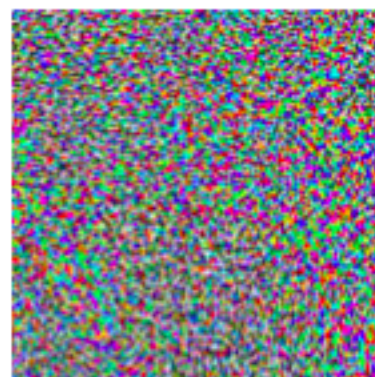**15 March 2022**

# Adversarial Machine Learning is an active area, and no winners



$+ .007 \times$

$=$

$\boldsymbol{x}$

"panda"
57.7% confidence

$\mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"nematode"
8.2% confidence

$\boldsymbol{x} +$
$\epsilon\,\mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"gibbon"
99.3 % confidence

# AML of Safeguards Images

Trained Model Results

Adversarial Results

"Hyperboloid Cooling Tower"

"NOT Hyperboloid Cooling Tower"



Carlini & Wagner* Attack



The trained model WAS NOT changed – rather, a small, nearly imperceptible, perturbation was added to the original image

"NOT Hyperboloid Cooling Tower"

"Hyperboloid Cooling Tower"





* N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *IEEE Symposium on Security and Privacy*, San Jose, 2017.

Error

# Frames of an image can also be fooled
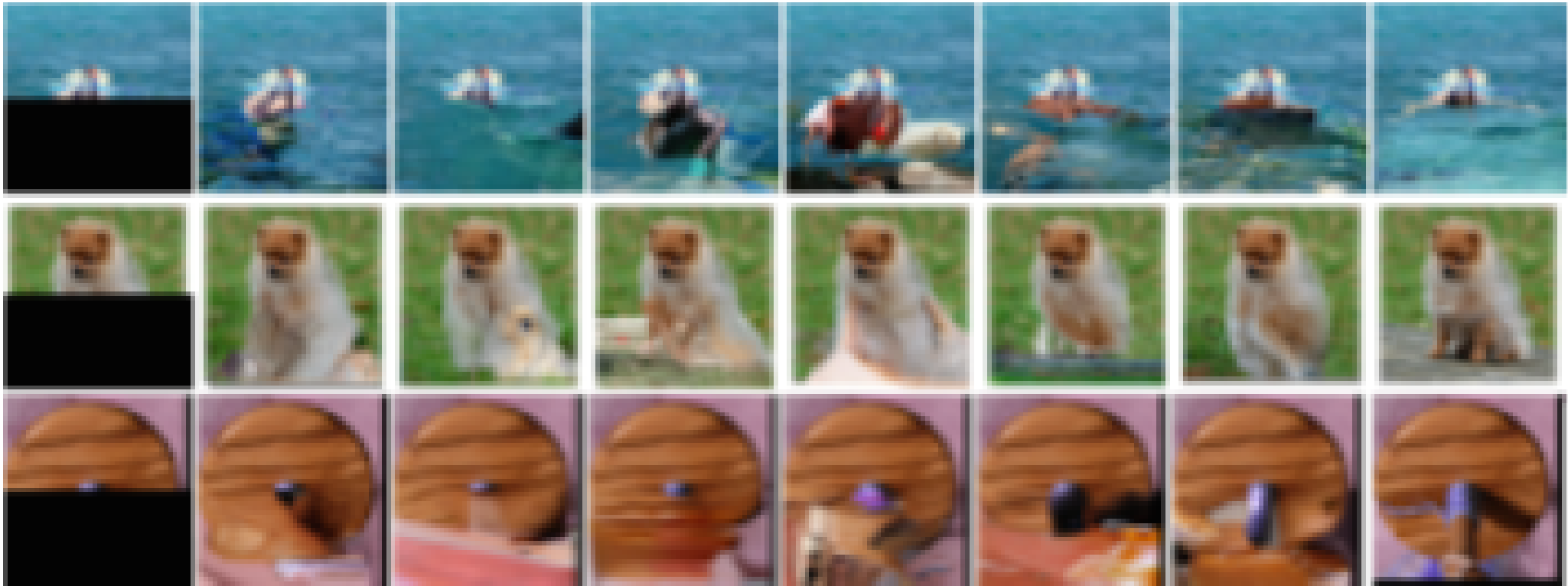
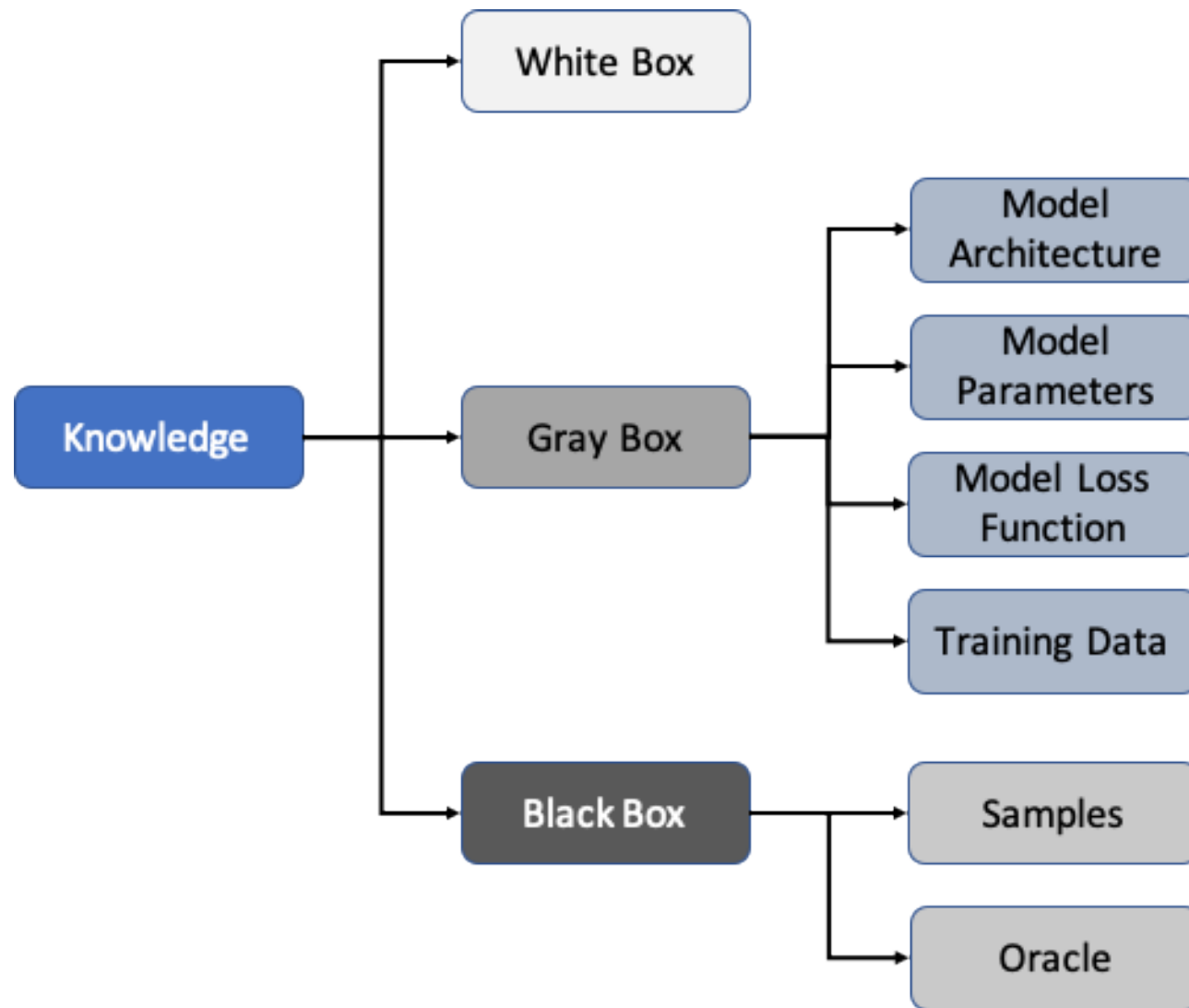# Can complete a partial image!
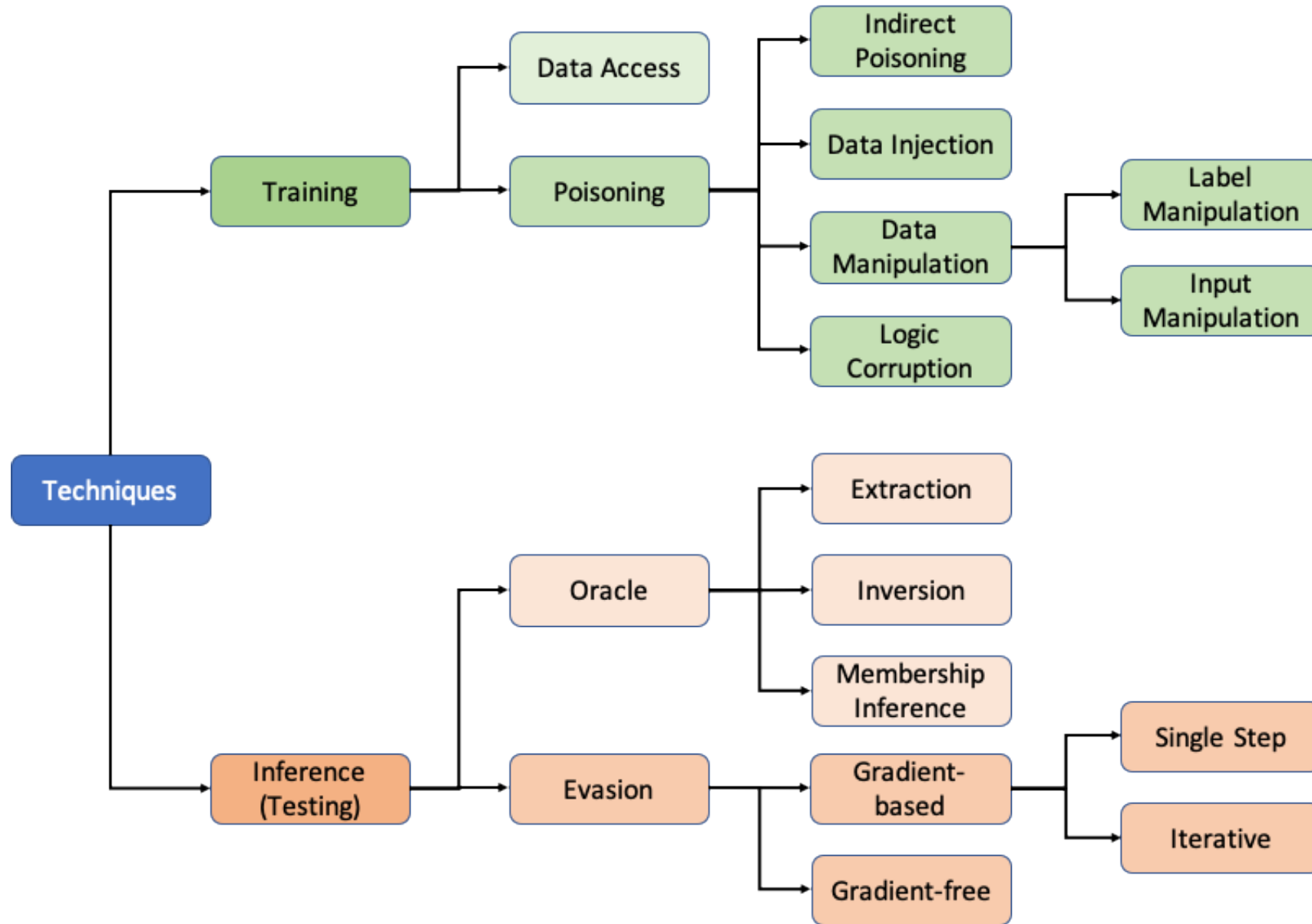


occluded        completions        original

# What is the adversary's knowledge?

# How does to adversary conduct an attack?

# Successful Adversarial Attacks can take several forms