



# HEMLOCK

POISON YOUR IMAGES

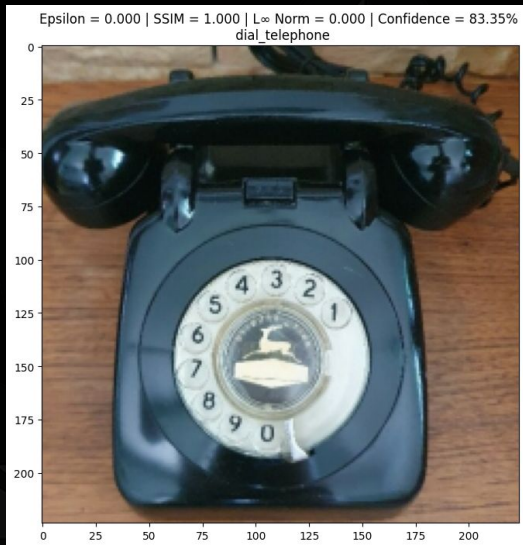
# HEMLOCK: THE FIGHT AGAINST AI

Hemlock is a concept of an interface that allows users to automatically apply poisoning techniques to their content. This app allows social media influencers to protect their work from being replicated or used by big tech companies to train AI models

# OVERALL APPROACH

- Decided to use the existing MobileNetV2 image classification model, & TensorFlow and Keras libraries.
- Experimented with different epsilon values to evaluate the model.
- Adversarial Attack Methods: Fast Gradient Sign Method, Projected Gradient Descent Method, and Carlini & Wagner Approach.

# IMAGES BEFORE AND AFTER

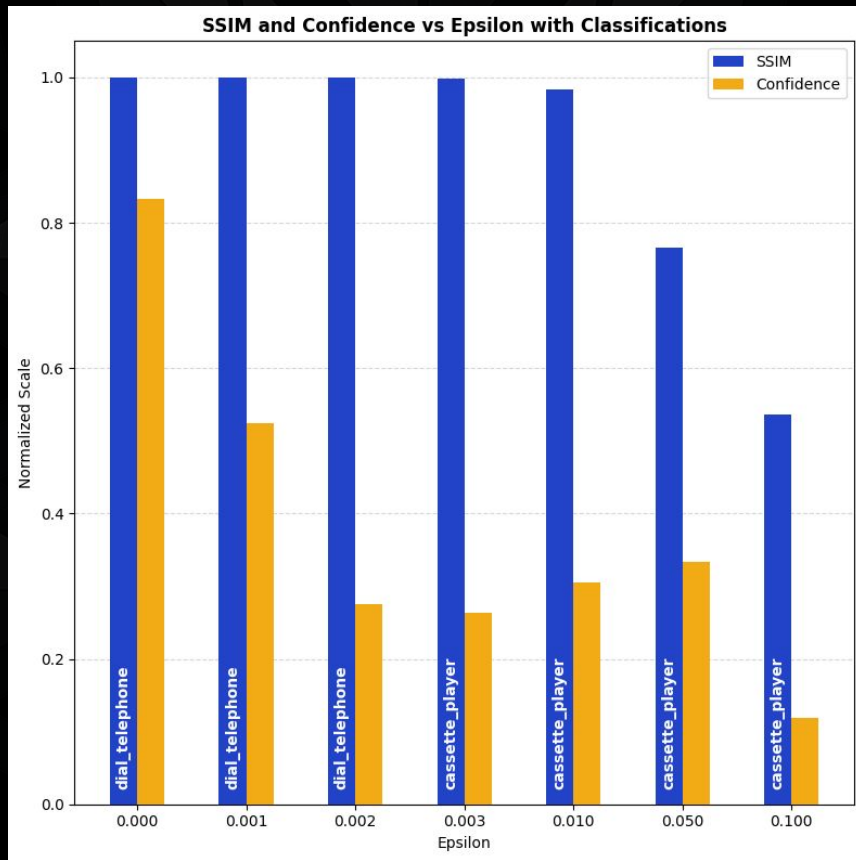


dial\_telephone + perturbations = cassette\_player

# PERFORMANCE METRICS

Structural Similarity Index (SSIM) is a comparison of image structure, brightness and contrast.

Confidence shows how sure the AI is about the identification.





DEMO

```

5
6 # import dependencies
7 import tensorflow as tf
8 import matplotlib as mpl
9 import matplotlib.pyplot as plt
10 from pathlib import Path
11 from PIL import Image
12 from rembg import remove
13 import numpy as np
14 import io
15 import tensorflow_hub as hub
16
17 #####
18 # # Helper function to preprocess the image so that it can be inputted in MobileNetV2
19 # def preprocess(image):
20 #     image = tf.cast(image, tf.float32)
21 #     image = tf.image.resize(image, (224, 224), method=tf.image.ResizeMethod.AREA)
22 #     image = tf.keras.applications.mobilenet_v2.preprocess_input(image)
23 #     image = image[None, ...]
24 #     return image
25
26 #####
27 # # Helper function to extract labels from probability vector
28 # def get_imagenet_label(probs):
29 #     return decode_predictions(probs, top=1)[0][0]
30
31 #####
32 # Helper function to remove the background from an image tensor
33 def remove_background(image, background_color=(255, 255, 255)):
34     """
35     Removes the background from a TensorFlow image tensor, crops to a square bounding box,
36     and replaces the background with a specified color.
37
38     Args:
39         image (tf.Tensor): Input image tensor.
40         background_color (tuple): RGB color for the new background (default is white).
41
42     Returns:

```

# PHASE 2

1. Add more loss functions to the model
2. Experiment with other attack methods
3. Calculate Attack Success Rate
4. Create UI



Q&A

# APPENDIX

[https://github.com/Jeff-Oliver/hemlock\\_image\\_poisoner.git](https://github.com/Jeff-Oliver/hemlock_image_poisoner.git)