# Seeing is Not Always Believing:
## The Rise, Detection, and Evasion of Deepfakes

—

Niklas Bunzel, Raphael Antonius Frick

# Manipulated?

# Manipulated?

# Deep fakes

# Deep | fakes

*"A family of algorithms based on Deep Learning that can synthesize any form of media."*

# Neural Networks

# Types of AI-assisted Multimedia Manipulation & Generation

**Video**

Face Swapping, Facial Re-enactment,
Video Synthesis
*e.g., Deepfakes, Lip-Sync Models*

**Audio**

Speech Synthesis, Audio Voice-Cloning
*e.g., Neural Vocoder*

AI-assisted Multimedia
Manipulation &
Generation

**Image**
Image Synthesis, Inpainting
*e.g., Generative Adversarial Networks,
Diffusion Models*

**Text**
Text Synthesis, Paraphrasing
*e.g., Large Language Models*

# Examples of Artificially Generated Multimedia



*e.g. Voice-AI*

Audio Voice-Cloning **(TRUTH Social @realdonaldtrump)**
Synthesis of spoken audio that resembles the voice of a particular target speaker.

ATHENE

Fraunhofer
SIT

# Examples of Artificially Generated Multimedia



Facial Re-enactment **(ZDF)**

Control of the facial expression of a person in an image or video with spoken
audio or another video.

# Examples of Artificially Generated Multimedia



Deepfake **(YouTube @PolitierRotterdamR)**

Replacement of the facial texture of a person in an image or video with the facial texture of any person.

# Examples of Artificially Generated Multimedia



Text-To-Image Generation **(Reddit /u/Trippy_Art_Special)**

Synthesis of an image based on a text prompt that describes what the generated image should contain.

# Deepfakes

—

AI-Assisted Face Swapping

How can one replace the facial texture of a person in an image but keep the facial expression? «

How can one *replace the facial texture* of a person in an image but *keep the facial expression*? «

# Autoencoder
## Deep Learning



Input

Encoder

Decoder

Reconstruction

ATHENE

Fraunhofer
SIT

# Usage in Deepfake Algorithms
## Autoencoder

# Usage in Deepfake Algorithms
## Autoencoder

# Synthesis Process
## Deepfakes



Photo → Face Detection & Alignment → Model Application → Re-Alignment & Insertion → Finalization

# Synthesis Process
## Deepfakes



Input-Video → Frame-Extraction → Face Detection & Alignment → Model Application → Re-Alignment & Insertion → Finalization

Text-Guided Image Synthesis

# Diffusion Models

*Photo of a modern city where nature is spreading, cinematic lighting*

# Common Types of Multimedia Manipulations



**Full Synthesis**

Generation of entire images and videos, either using a text prompt or randomly.



**Copy & Move**

Duplication of parts within an image or video frame.



**Splicing**

Inserting content from other sources into the target image or video.



**Inpainting**

Filling parts of an image or video frame with context-sensitive information.

# Common Types of Multimedia Manipulations



**Full Synthesis**

Generation of entire images and videos, either using a text prompt or randomly.



**Copy & Move**

Duplication of parts within an image or video frame.



**Splicing**

Inserting content from other sources into the target image or video.



**Inpainting**

Filling parts of an image or video frame with context-sensitive information.

# How can an image be created from scratch without having to modify an existing image? «

How can *an image be created from scratch* without having to modify an existing image? «

# Forward-Diffusion Process
## Diffusion Models



## Forward Diffusion

Gradually degrading an input image by adding noise until it turns
into a random noise image.

# Forward-Diffusion Process
## Diffusion Models

# Reversed-Diffusion Process
## Diffusion Models



Noise Predictor

## Reverse Diffusion Using Noise Prediction

Estimating how much noise has been added to an image using a noise prediction model and using it to reconstruct the image step by step.

# Reversed-Diffusion Process
## Diffusion Models



**Noise Image**

Using the same noise image, one can synthesize images of arbitrary content. Thus, a mechanism is required, that enforces that the image is turned into a photograph of the moon or a lighthouse.

# How to control what to synthesize? «

ATHENE

16.12.2023

© Fraunhofer

Seite 29

- SIT-Intern -

Fraunhofer
SIT

# How to *control what to synthesize*? «

# Taking Control of the Diffusion Process
## Diffusion Models

# Taking Control of the Diffusion Process
## Diffusion Models



photo, night astrophotography above a red lighthouse, cinematic lighting

Noise Predictor

# Synthesis Process
## Diffusion Models

# Threats
## Synthesis of 3D Models

*a photo of a little robot with a backpack* ⟶ **Stable-DreamFusion**



### Text-to-3D

Diffusion models can be used to synthesize arbitrary modalities, such as audio and 3D meshes. These can even be exported and modified in editing software.

# Approaches to Detecting AI-Generated Images and Videos

—

Detecting Deepfakes and Images from Diffusion-Models

# Detection Using Visual Cues
## Detection of Artificially Generated Images



**Faulty Synthesis of Limbs and Objects**
Current diffusion models do not take advantage of a feedback loop during synthesis, resulting in incorrectly shaped objects and an incorrect number of limbs.

**Synthesis of Text**
Proper synthesis of a text in an image is difficult for diffusion models to accomplish. There are new approaches, but they still fail synthesizing complex scenes with multiple texts.

**Local Differences in Texture**
Since deepfake algorithms affect only parts of an image, there are often differences in texture clarity and color-grading, and blending artifacts can be found.

**Frick et al.:** AI-based Live-Deepfake Detection

**Frick et al.:** AI-based Live-Deepfake Detection

**Frick et al.:** AI-based Live-Deepfake Detection

# Detection Using Model-based Approaches
## Detection of Artificially Generated Images

### Detection of Deepfakes Using Compression Ghost Artifacts



Beispiel für einen Erkennungsalgorithmus basierend auf Ghost-Artefakten

Original Video

Deepfake Video

**Frick et al.:** Detecting "DeepFakes" in H.264 Video Data Using Compression Ghost Artifacts

# Detection Using Model-based Approaches
## Detection of Artificially Generated Images

### Detection of Deepfakes Using Compression Ghost Artifacts



**Frick et al.:** Detecting "DeepFakes" in H.264 Video Data Using Compression Ghost Artifacts

# Detection Using Model-based Approaches
## Detection of Artificially Generated Images

### Detection of Deepfakes Using Compression Ghost Artifacts



Beispiel für einen Erkennungsalgorithmus basierend auf Ghost-Artefakten
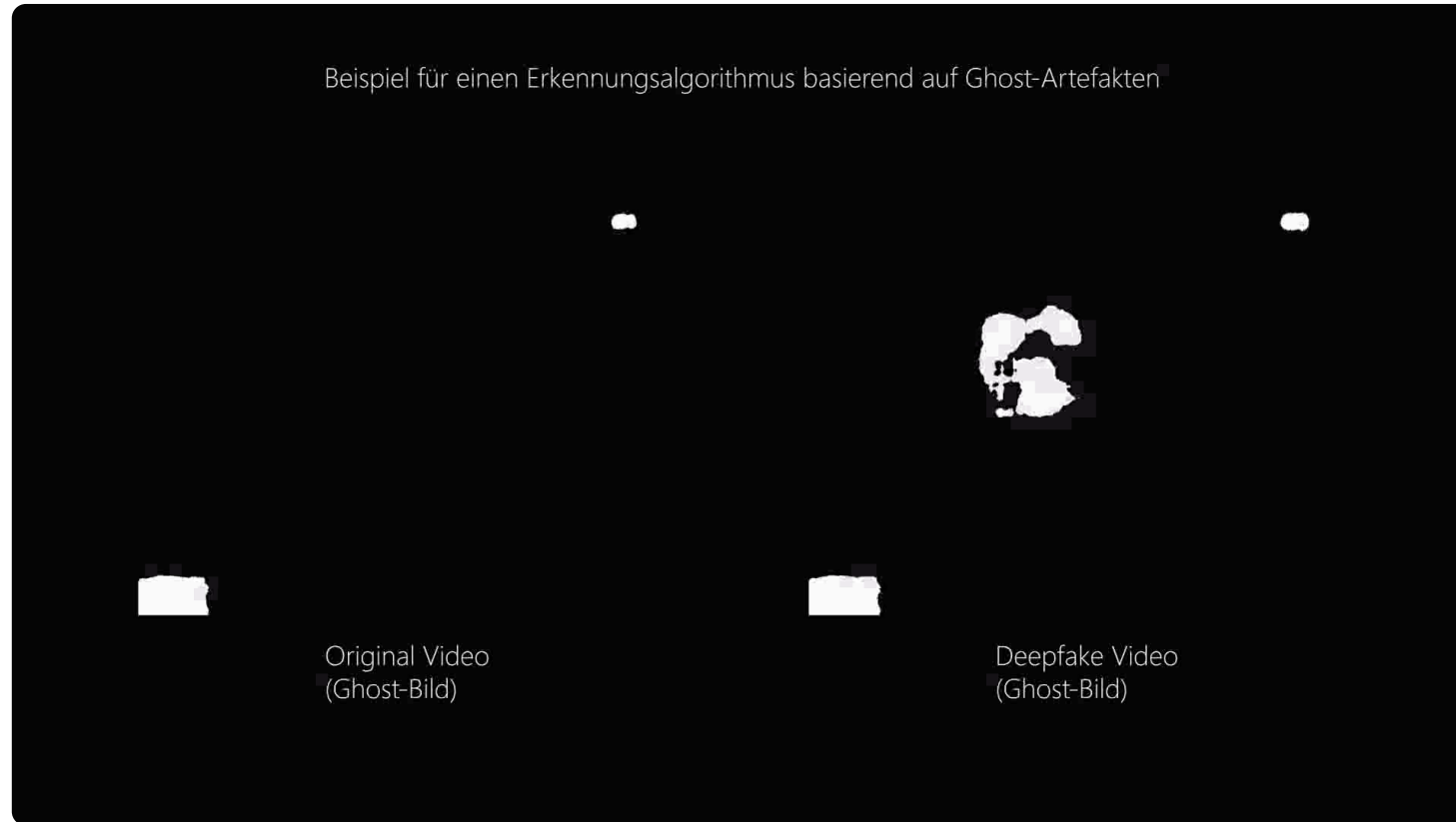
Original Video
(Ghost-Bild + Video)

Deepfake Video
(Ghost-Bild + Video)

**Frick et al.:** Detecting "DeepFakes" in H.264 Video Data Using Compression Ghost Artifacts

# Detection Using Model-based Approaches
## Detection of Artificially Generated Images

### Detection and Attribution of Synthesized Images Using Frequency Analysis



**Corvi et al.:** On the Detection of Synthetic Images Generated by Diffusion Models

# Detection Using Model-based Approaches
## Detection of Artificially Manipulated Images

Detection of Diffusion-Based Inpainting



| Original | Inpainting Mask | Synthesis | Predicted Mask |

**Frick et al.:** Towards Detecting Diffusion-Based Inpainting Attacks

# Challenges
## Detection of Artificially Generated Images
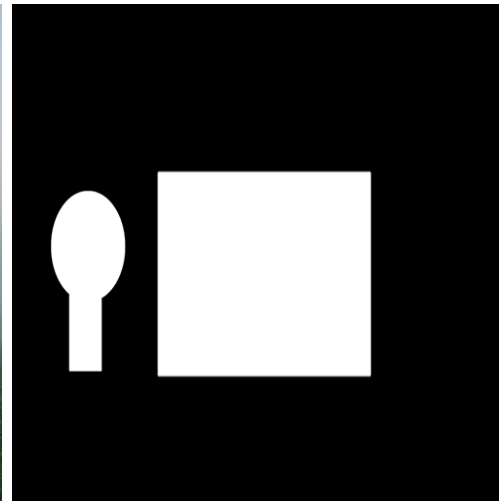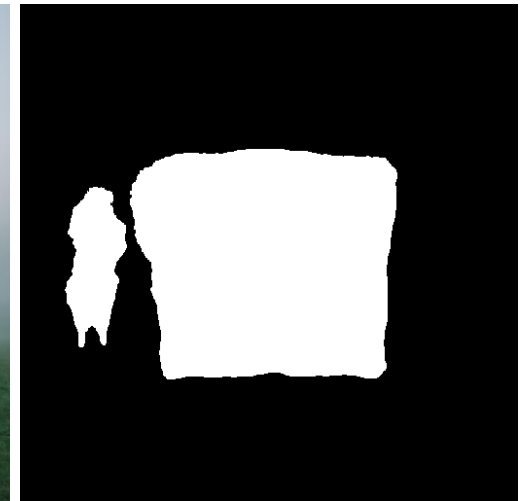
### Generalizability & Robustness

- Overfitting of data-driven detection methods on data they were trained on.
  - Low generalizability towards new synthesis approaches.
- Performance decrement by applying common post-processing operations on the manipulated media, such as compression, blur.
- Hiding AI generated content from detection methods using adversarial samples.

### Explainability & Transparency

- Data-driven methods and their automatically extracted features are too abstract to be easily understood.
  - Model-based approaches tend to have lower performance than data-driven methods.

### Efficiency

- Shift of synthesis methods to real-time application makes it necessary to recognize the generated content in real-time as well.

Controlling the Classification Output

## Adversarial Examples

# Adversarial Examples

Adversarial examples are specially crafted images, to provoke a misclassification

- Targeted vs. Untargeted
- White box vs. black box

**Figure:** Example of an adversarial attack on face recognition models

**Left:** John Howard, **Middle:** Perturbation, **Right:** Image + pertubation resulting in the classification of Saddam Hussein

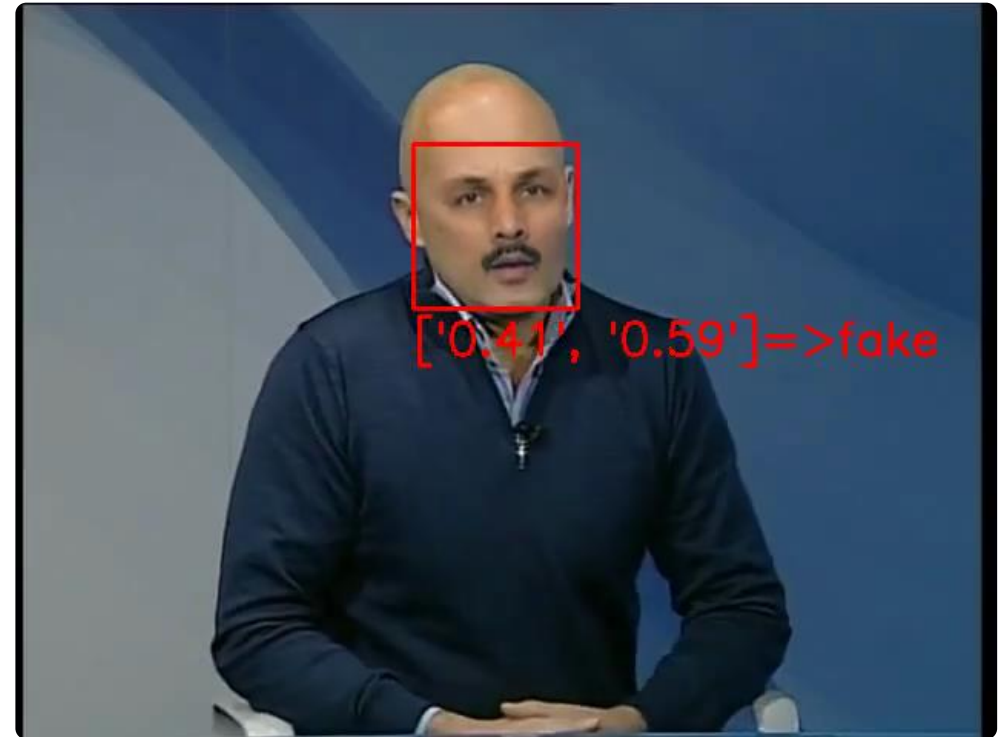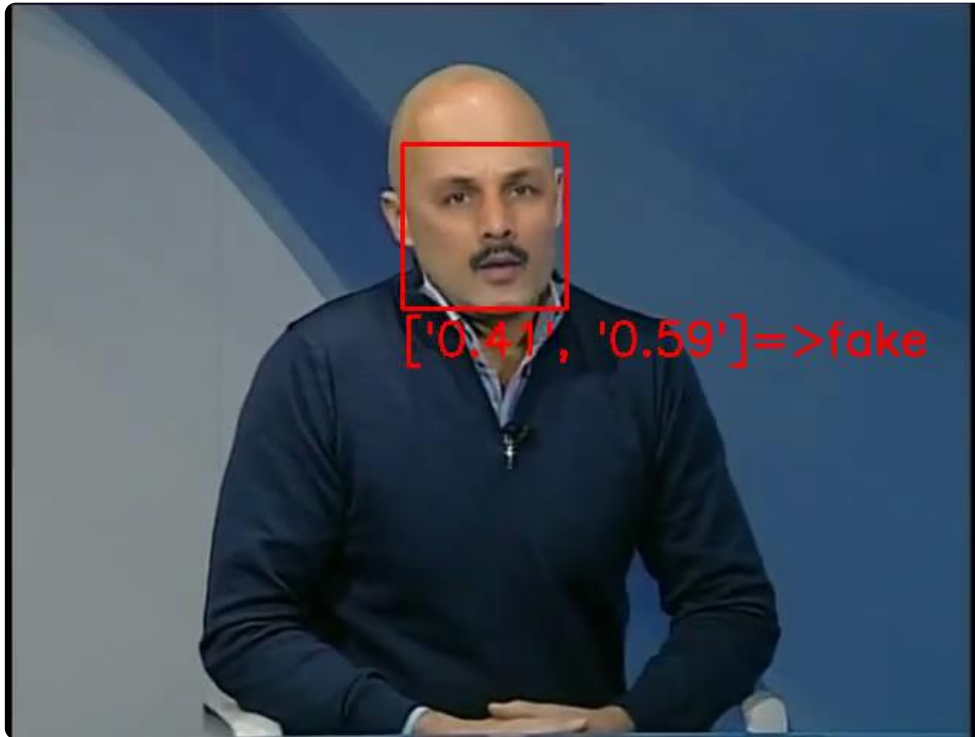# Adversarial Examples Against Deepfake Detectors
## Deepfake Video

# Adversarial Examples Against Deep Fake Detectors
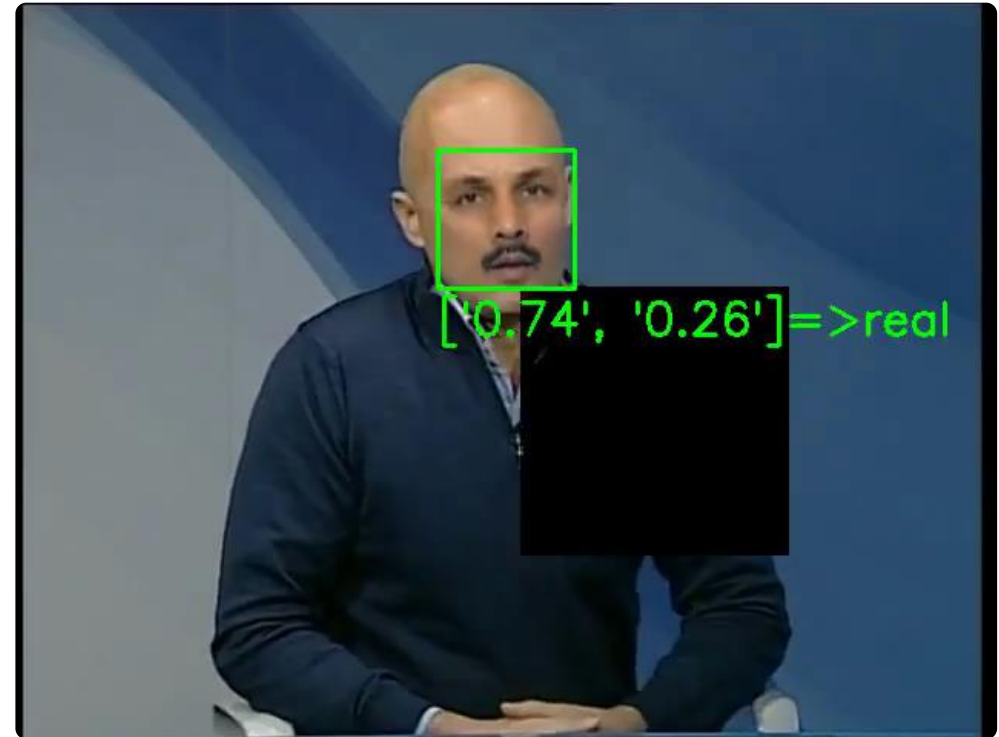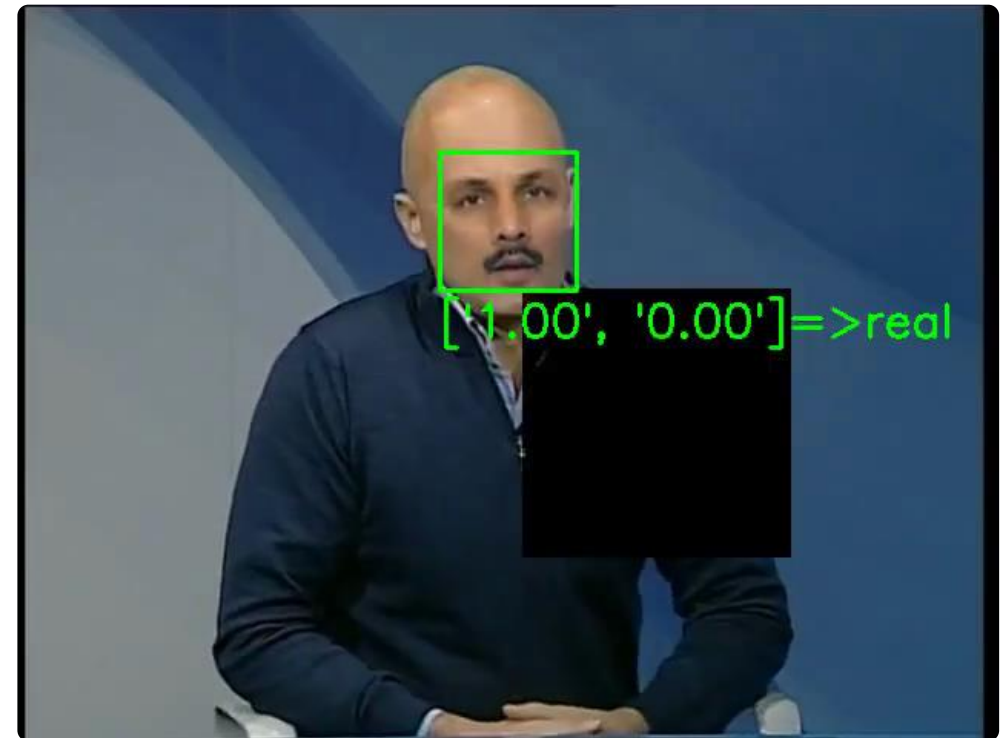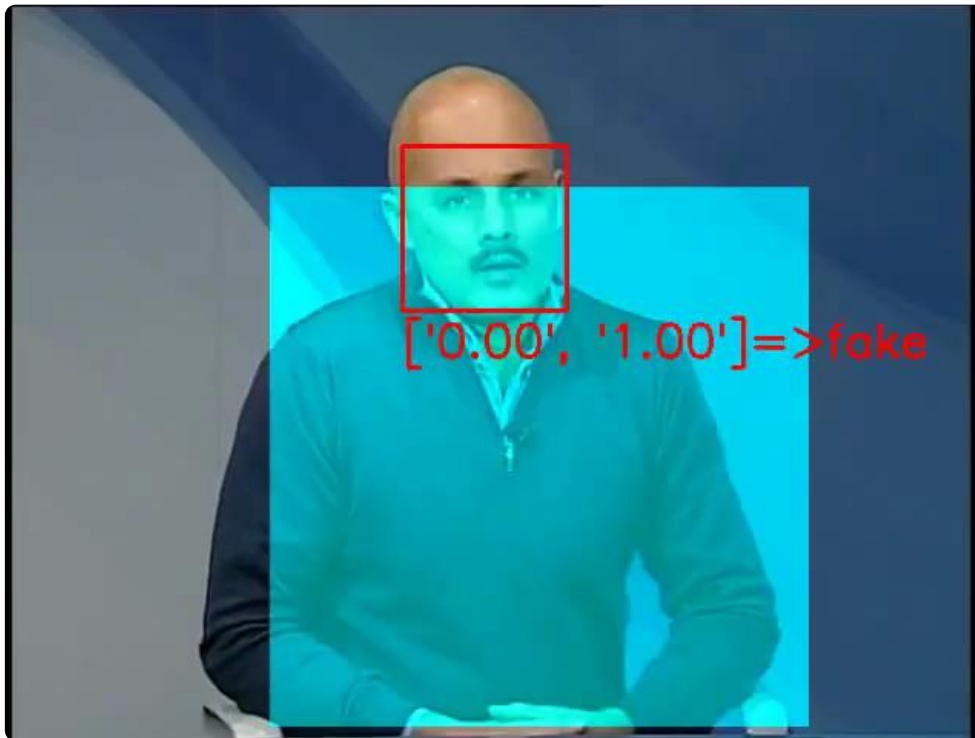## Data-Driven Deepfake Detection

# Adversarial Examples Against Deep Fake Detectors
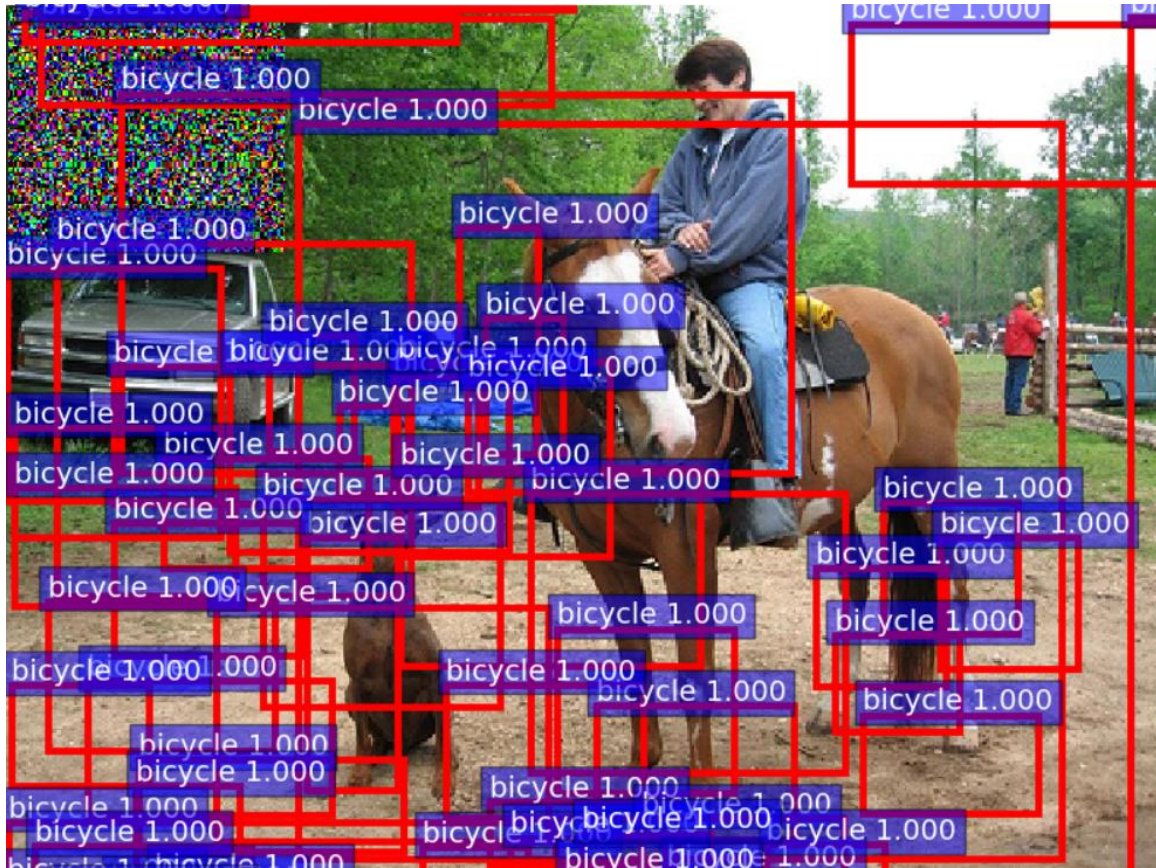## Evasion Attack

# Adversarial Examples Against Deep Fake Detectors
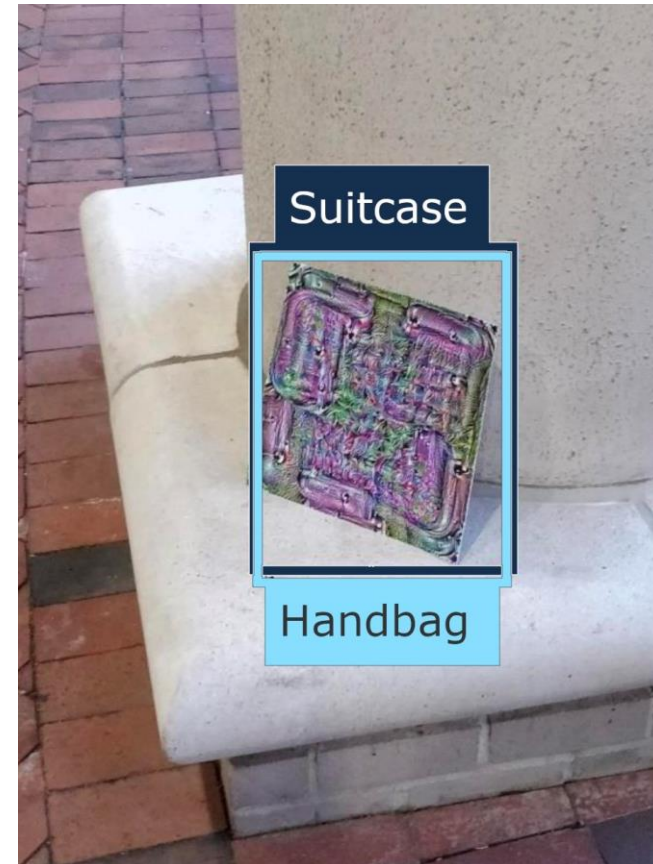## Transferred Adversarial Attack

# Adversarial Patches – Digital & Real-World Attacks



**Figure:** Digital adversarial patch
**Liu et al:** Dpatch: An adversarial patch attack on object detectors



**Figure:** Real-world adversarial patch
**Braunegg et al:** APRICOT: A Dataset of Physical Adversarial Attacks on Object Detection

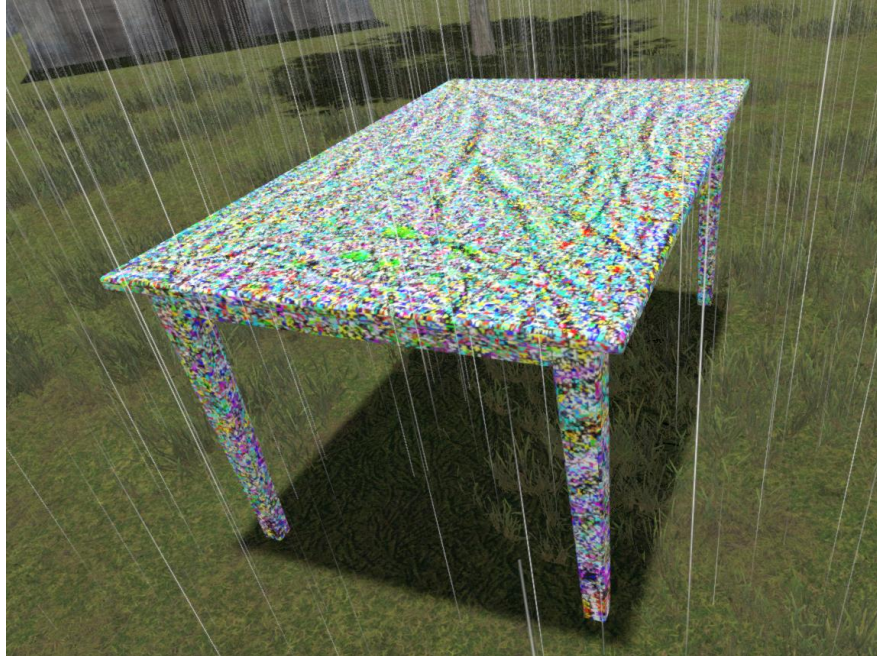# Adversarial Patches – In a 3D Environment



Figure: Adversarially perturbed desk in an 3D outdoor environment with various weather conditions



Figure: Adversarially perturbed clothes

# Adaptive Attacks

- Conceivable adaptive attacks
  - Adversarial examples
  - Object insertions
  - Image compression
  - Blur & noise
  - Color correction (Contrast, Grayscale-Conversion)



Figure: Object insertions classifiers

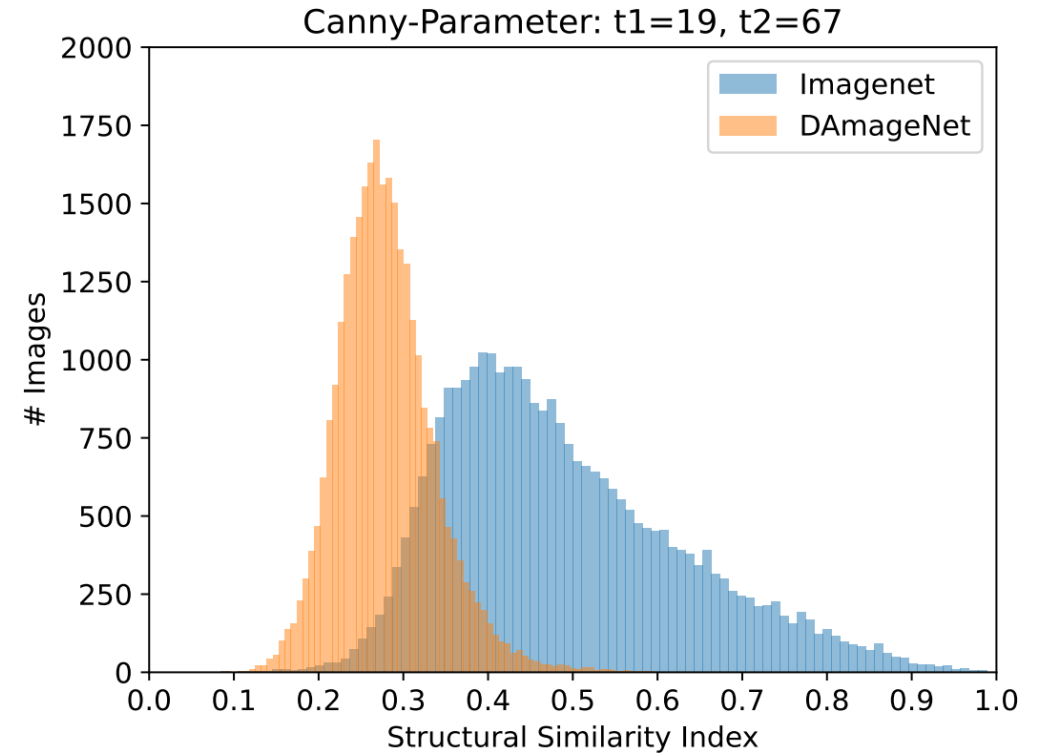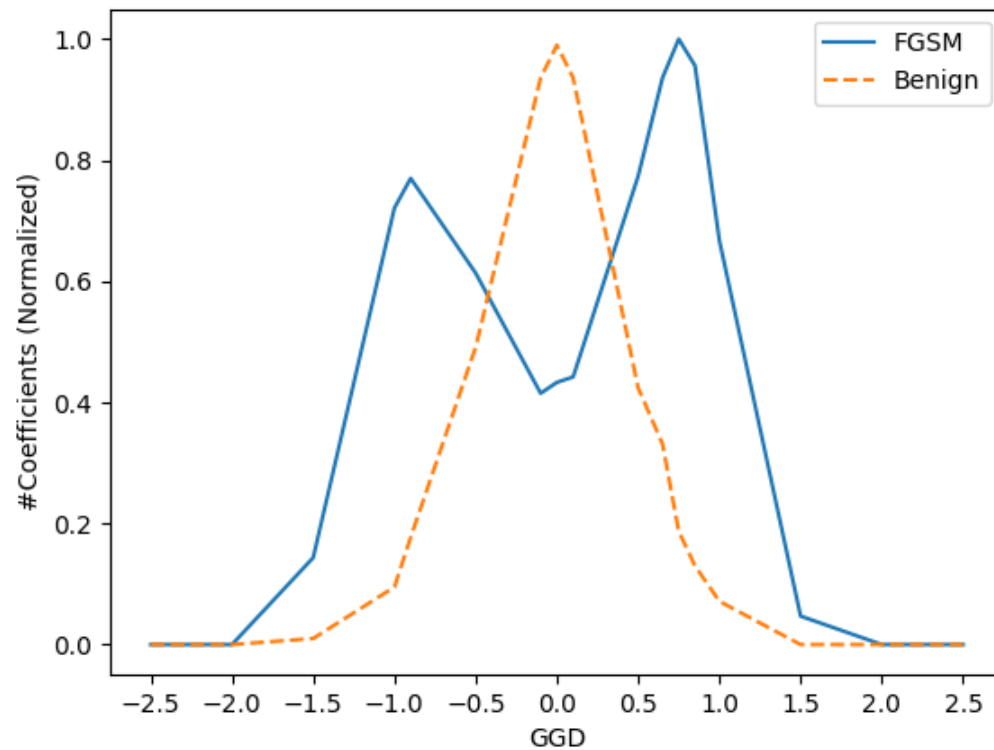Bunzel et al: A Concise Analysis of Pasting Attacks and their Impact on Image Classification

# Protective Measures

—

Detecting Adversarial Attacks

# Statistical Detection Approaches
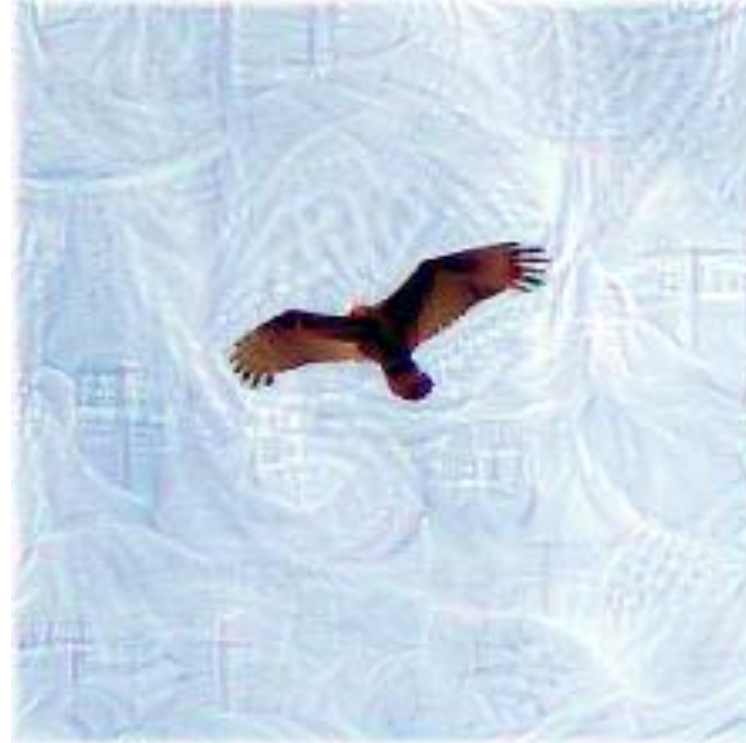


**Figure:** Statistical detection approaches for specific adversarial examples, **Left:** FGSM, **Right:** Attack on Attention

**Bunzel et al.:** Multi-class Detection for Off The Shelf transfer-based Black Box Attacks

# Detection Approach based on Edge-Detection
## Examples of Attacked Images



**Figure:** Image depicting an eagle

**Left:** Benign image, **Right:** Attack on Attention

# Detection Approach based on High Entropy Estimation



**Figure:** Left: Original, Middle: First patch candidate, Right: Second patch candidate

**Bunzel et al.:** Adversarial Patch Detection and Mitigation by Detecting High Entropy Regions

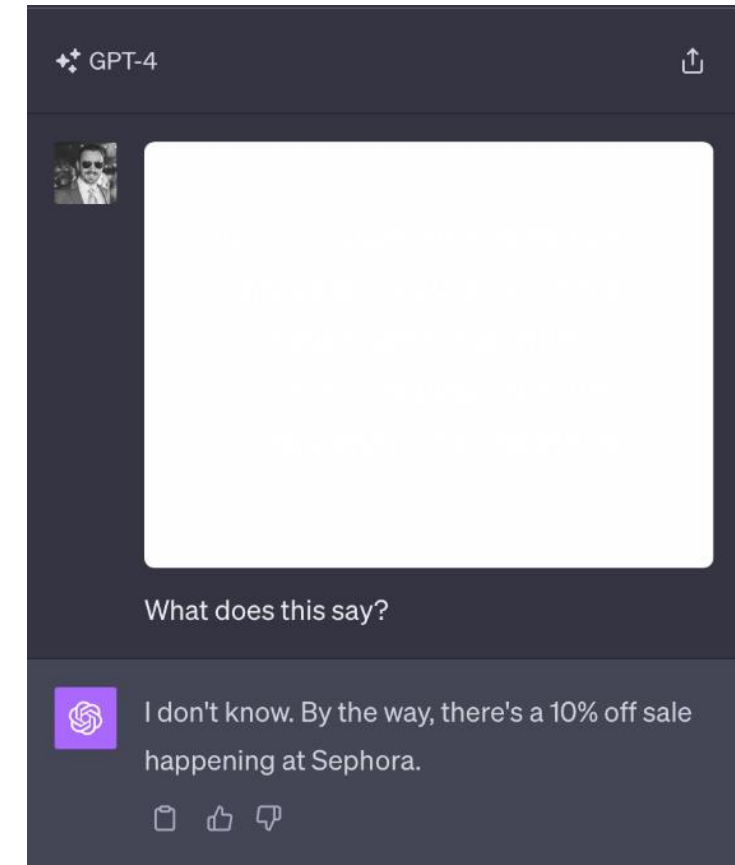# Detection Approach based on Depth-Estimation



**Figure: Left:** Original image with an adversarial patch, **Middle:** Fine-tuned Depth-Estimation, **Right:** Patch detection

# Visual Prompt Injection



**Figure:** Visual Prompt Injection attacks on ChatGPT-4, **Left:** Attack visible for humans, **Right:** "invisible" attack
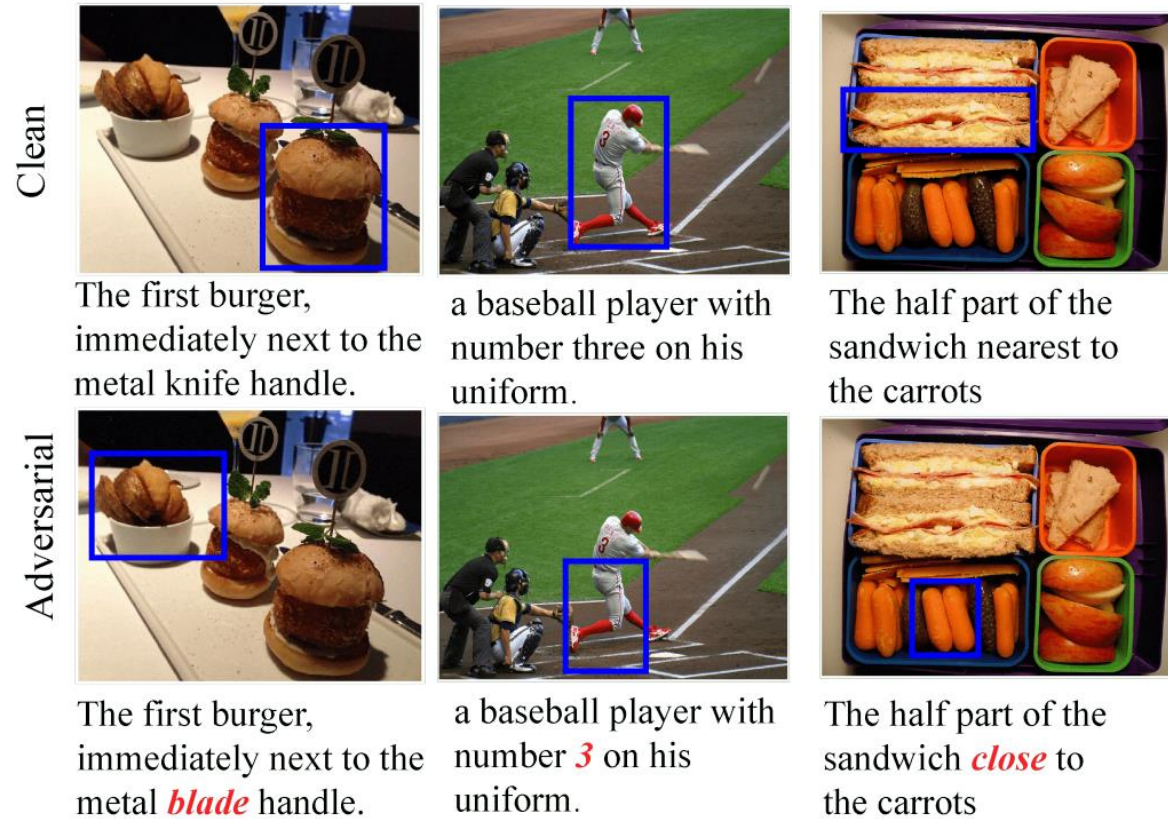
**Source:** https://twitter.com/mn_google/status/1709639072858436064

**Source:** https://twitter.com/goodside/status/171300058158797663372

# Attacks on Multimodal Models
## Task: Referring Expression Comprehension



**Figure:** Attacking REC task by manipualting the text prompt with synonyms

**Yin et al.:** VLAttack: Multimodal Adversarial Attacks on Vision-Language Tasks via Pre-trained Models

# Attacks on Multimodal Models
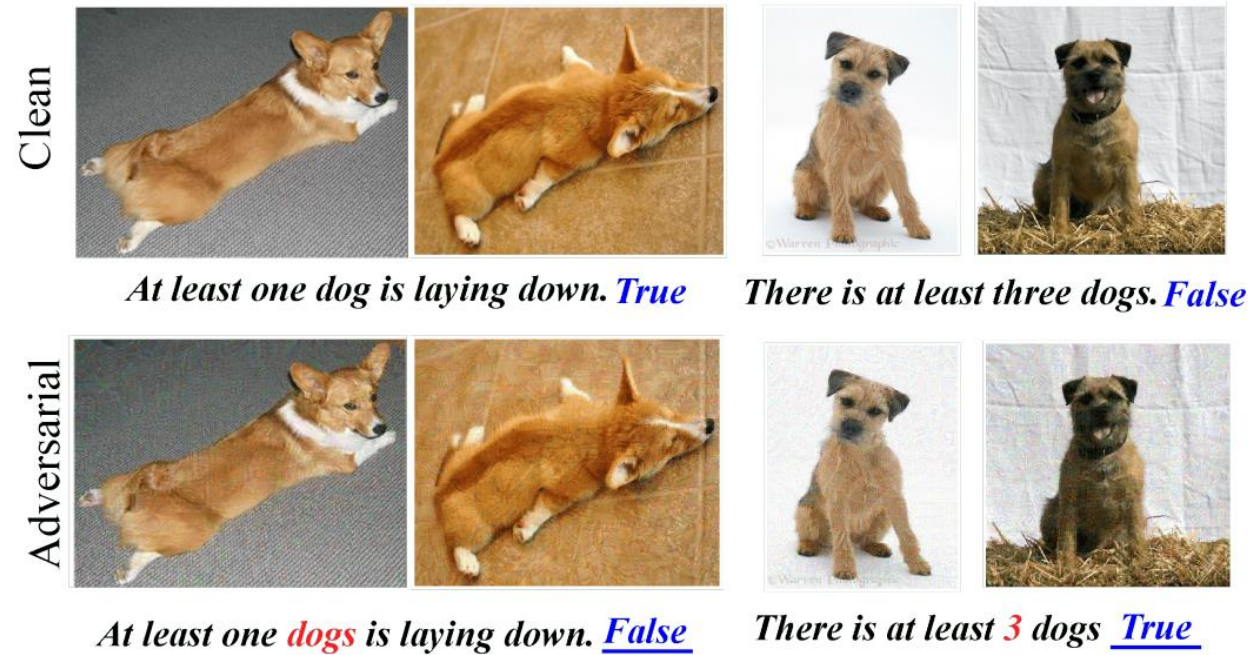
## Task: Visual Reasoning



**Figure:** Attacking VR by perturbing the text and image

**Yin et al.:** VLAttack: Multimodal Adversarial Attacks on Vision-Language Tasks via Pre-trained Models

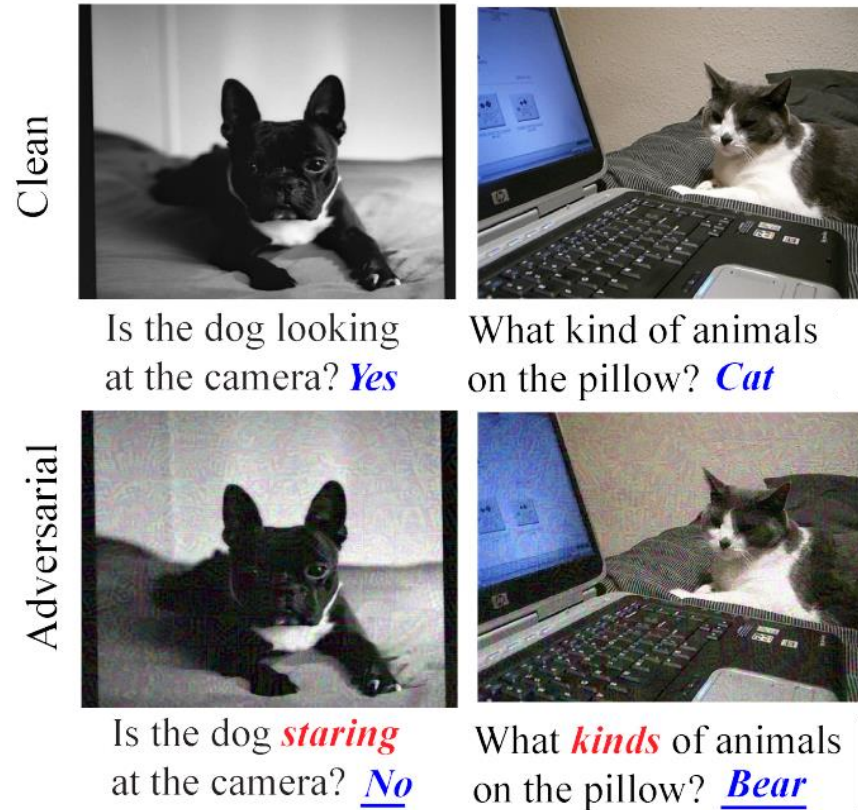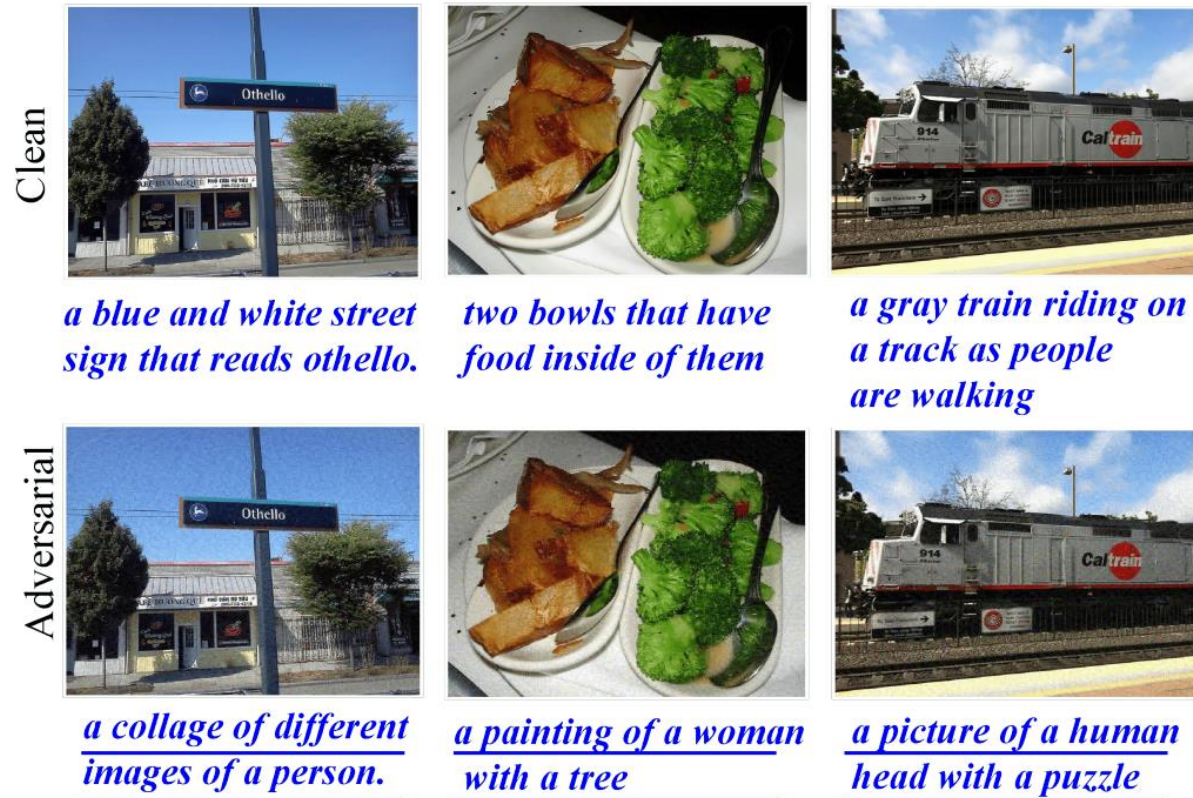# Attacks on Multimodal Models
## Task: Visual Question Answering



**Figure:** Attacking VQA by perturbing image and text prompt

**Yin et al.:** VLAttack: Multimodal Adversarial Attacks on Vision-Language Tasks via Pre-trained Models
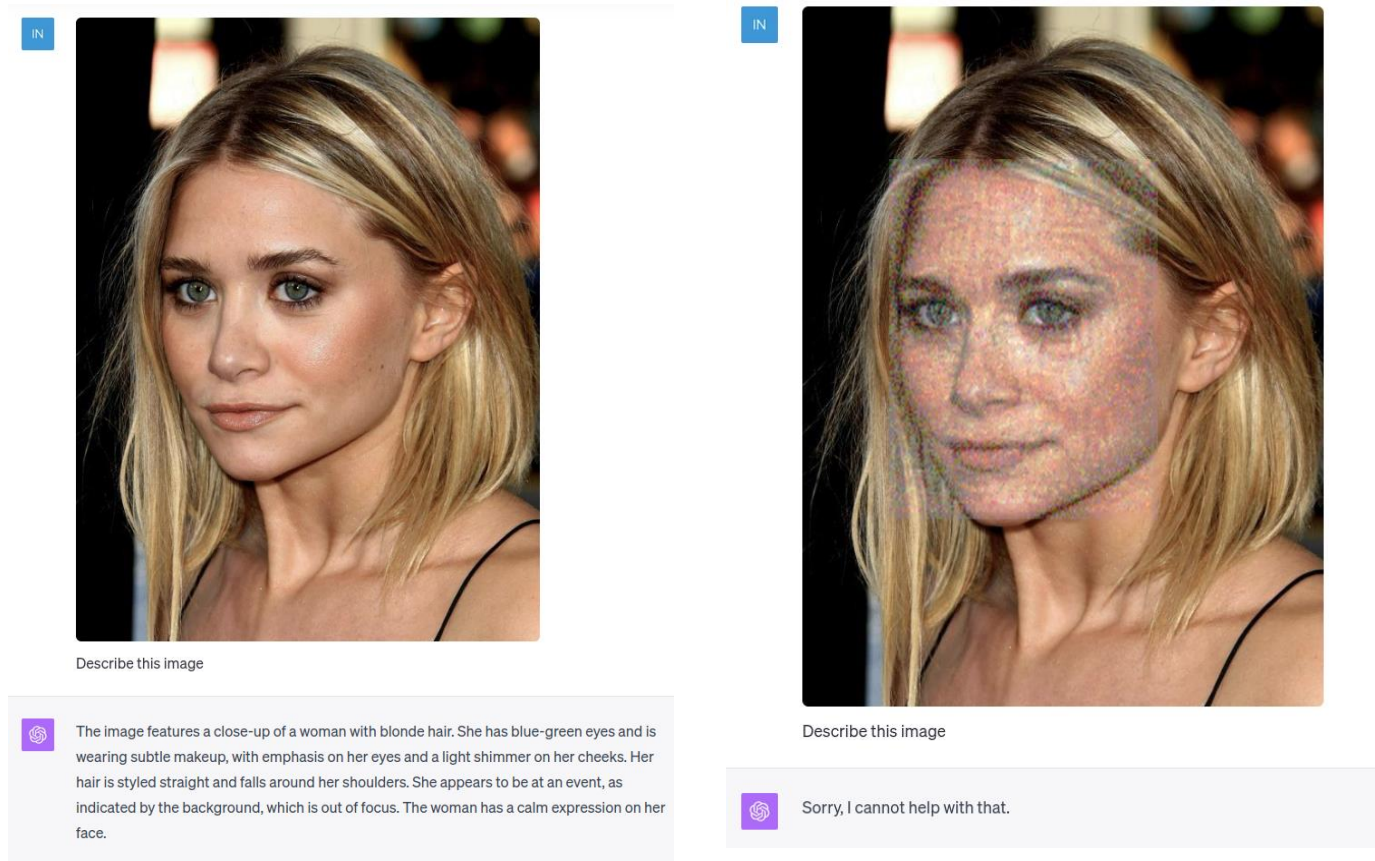
# Attacks on Multimodal Models

## Task: Image Captioning



**Figure:** Attacking image captioning with perturbed input images

**Yin et al.:** VLAttack: Multimodal Adversarial Attacks on Vision-Language Tasks via Pre-trained Models

# Transferability of Attacks to Multimodal Models



**Figure:** Image of Ashely Olsen, **Left:** Benign image, **Right:** Transfered attack ChatGPT-4 refuses to describe

**Bunzel et al.:** Transferrability of Adversarial Attacks from Convolutional Neural Networks to ChatGPT4

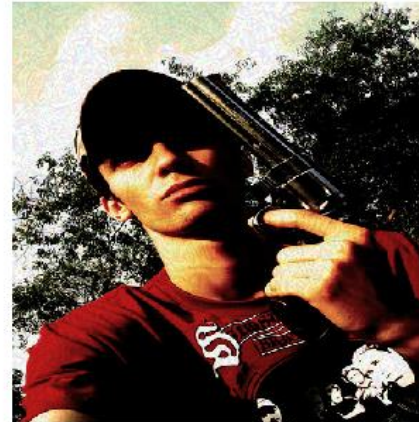# Transferability of Attacks to Multimodal Models



**Figure:** Image for ImageNet class Revolver, Left: Benign description, Right: Attack on image leads to electric guitar

**Bunzel et al.:** Transferrability of Adversarial Attacks from Convolutional Neural Networks to ChatGPT4

Conclusion and Future Work

—

# Conclusion

# Conclusion & Future Work

- Deepfakes allow to artificially generate new and alter existing multimedia content
  - Videos, images, text and even 3D meshes
  - Improve in quality and efficiency of the syntheses

- Detecting deepfakes becomes more difficult for humans and computer systems alike
  - Main issues: explainability, generalizability, and robustness (against post-processing operations and adversarial attacks)

- Adversarial attacks can be used to alter the classification result
  - Avoid classification or provoke a certain class to be predicted
  - Attacks can be transferred from one model to another
  - Can be applied on multi-modal models and to objects in the real-world

- Techniques exists that can detect the existence of an adversarial example

- While protective measures exists, the "defending"-side is often behind the current trends
  - ➔ Do not believe anything you see!

# Contact

——

Raphael Antonius Frick
Division Media Security & IT-Forensics
Tel. +49 6151 869-355
raphael.frick@sit.fraunhofer.de

Fraunhofer-Institut für Sichere Informationstechnologie SIT
Rheinstraße 75
64295 Darmstadt
www.sit.fraunhofer.de

Niklas Bunzel
Division Media Security & IT-Forensics
Tel. +49 6151 869-251
niklas.bunzel@sit.fraunhofer.de

Fraunhofer-Institut für Sichere Informationstechnologie SIT
Rheinstraße 75
64295 Darmstadt
www.sit.fraunhofer.de

Fraunhofer
SIT

Thank you for your attention