

# Spoofing Cybersecurity

Arming your organization against deepfakes and AI-synthesized data

**Adel Abusitta**









# Agenda

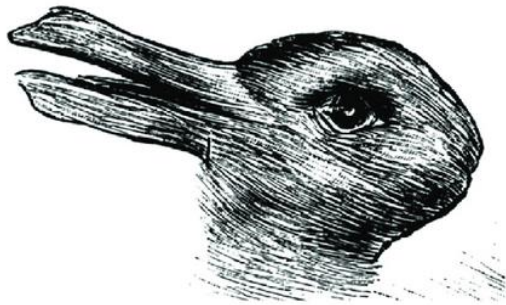
- The State of Deepfakes
- Generative Adversarial Networks and Deepfakes
- How to detect Deepfakes
- How to prepare an organization against Deepfakes

# Why does Deepfakes matter?

## Deepfakes matter because:

**Believability:** If we see and hear something with our own eyes and ears, we believe it to exist or to be true, even if it is unlikely.

- The brain's visual system can be targeted for misperception, in the same way optical illusions and bistable figures trick our brains.



Jastow rabbit-duck

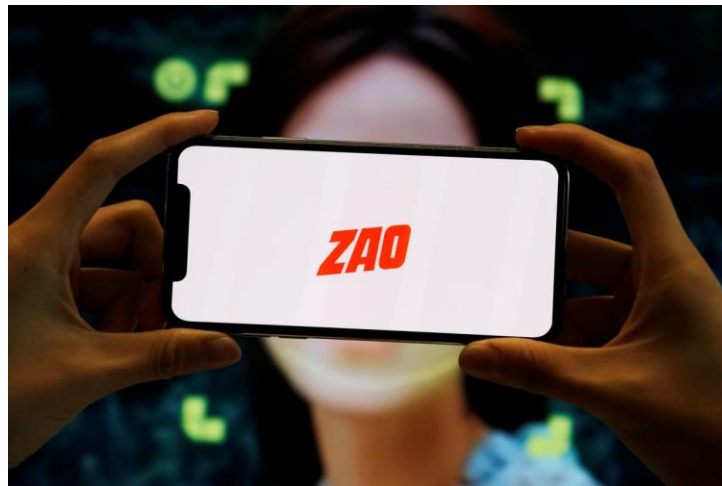


Rubin vase-faces

# Why does Deepfakes matter?

## Deepfakes matter because:

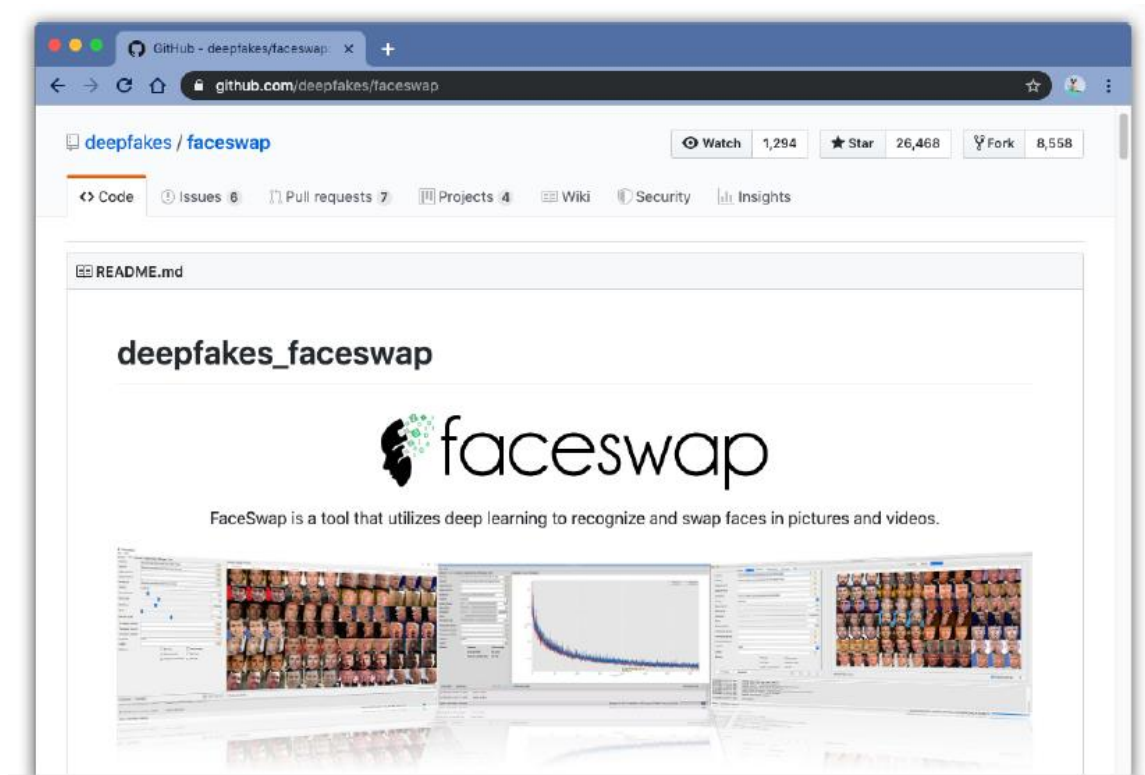
**Accessibility:** the technology of today and tomorrow, will allow all of us to create fakes that appear real, without a significant investment in training, data collection, hardware and software.



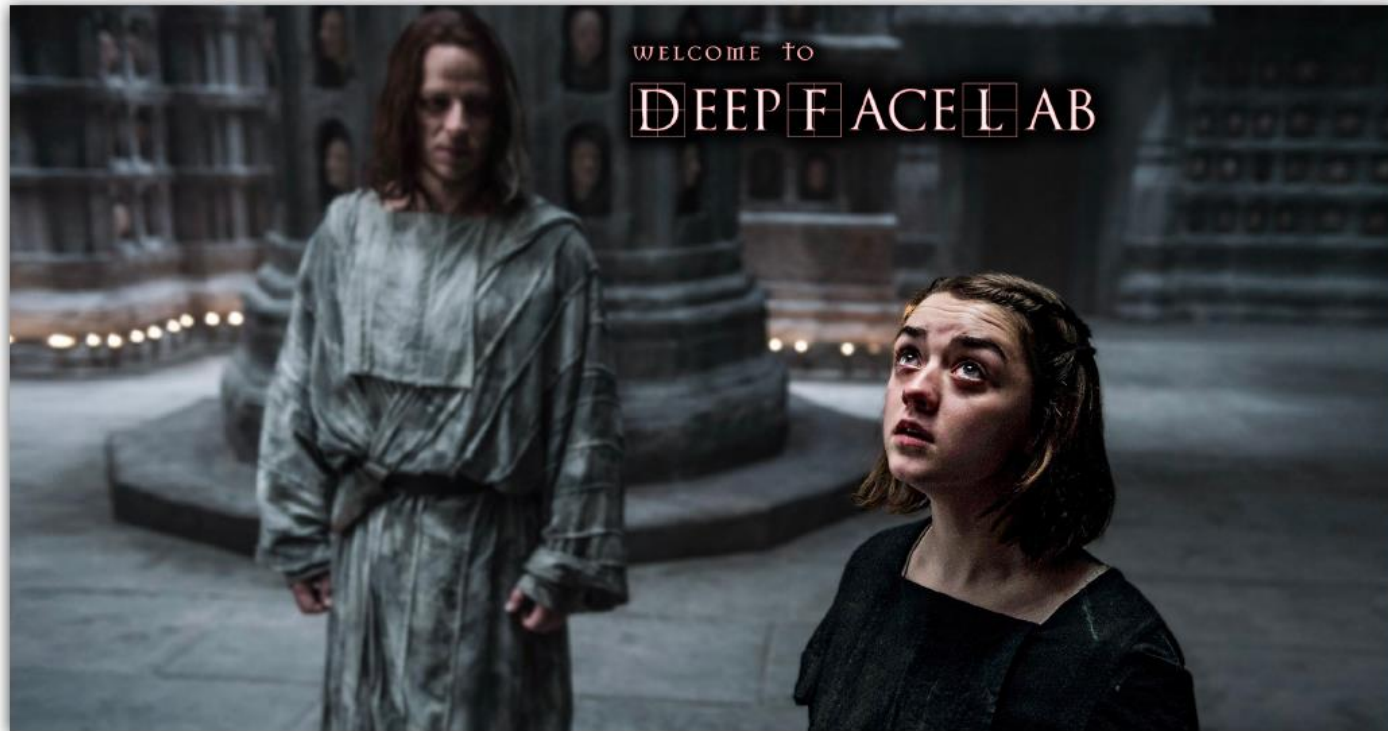


# Faceswap

Open Source multi-platform Deepfakes software.



# DeepFaceLab



# Zao,

The popular Chinese app for mobile devices lets users place their faces into scenes from movies and TV shows, for free.



# How do deepfakes work?

## How do deepfakes work?

Consider below deepfake featuring Jim Carrey and Alison Brie:



Original showing Alison Brie



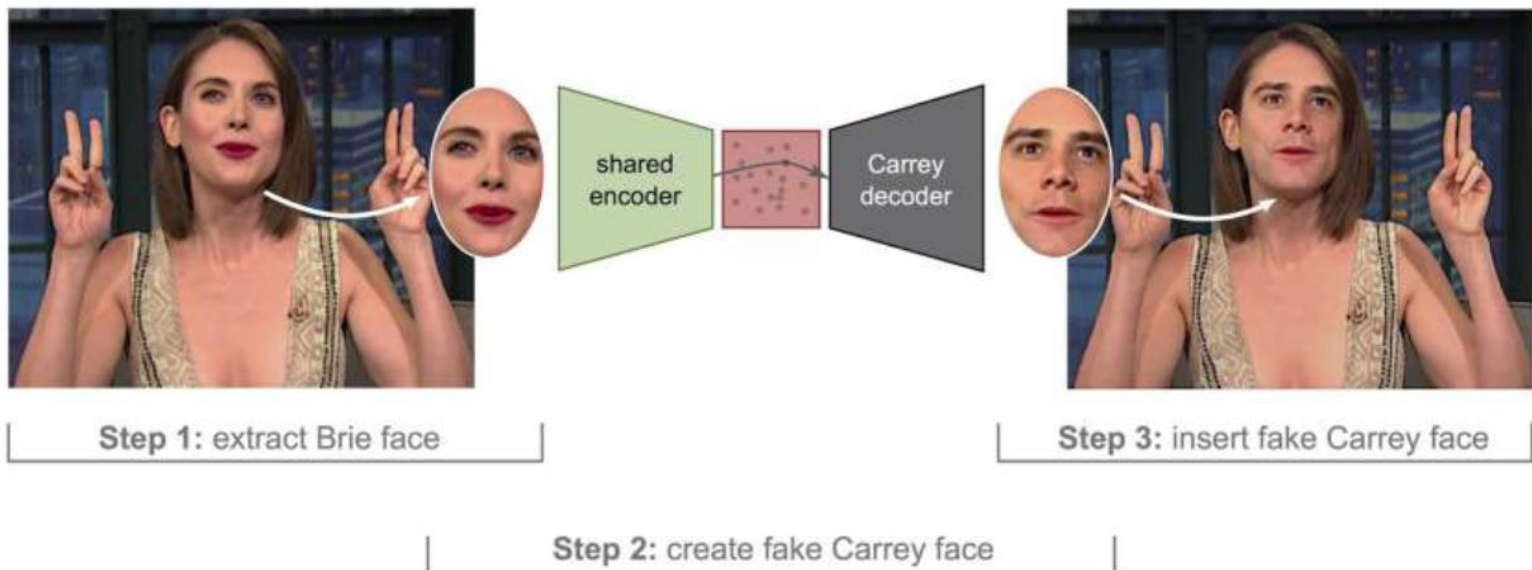
Deepfake showing Jim Carrey instead of Brie

The original Alison Brie video: <https://www.youtube.com/watch?v=QBmYDzLhWoY>

The deepfake with Jim Carrey: <https://www.youtube.com/watch?v=b5AWhh6MYCg>

# How do deepfakes work?

Many deepfakes are created by a three-step procedure:



# How do deepfakes work?

- **Step 1:** The image region showing Brie's face is **extracted** from an original frame of the video. This image is then used as input to a deep neural network (DNN)
- **Step 2:** The DNN automatically **generates** a matching image showing Carrey instead Brie
- **Step 3:** This generated face is **inserted** into the original reference image to create the deepfake

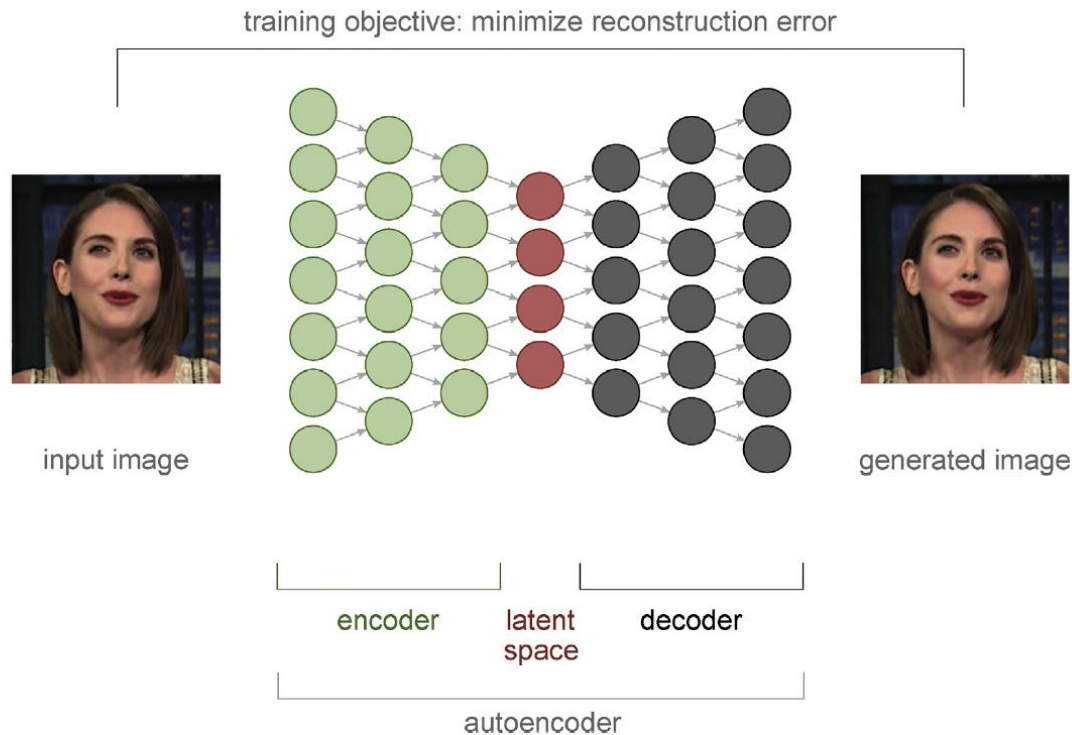
# How do deepfakes work?

- Deepfakes are commonly created using a specific deep network architecture known as **autoencoder**.
- Autoencoders consist of three subparts:
  - **an encoder** (recognizing key features of an input face)
  - **a latent space** (representing the face as a compressed version)
  - **a decoder** (reconstructing the input image with all detail)



# How do deepfakes work?

- **Autoencoder:** a DNN architecture commonly used for generating deepfakes.

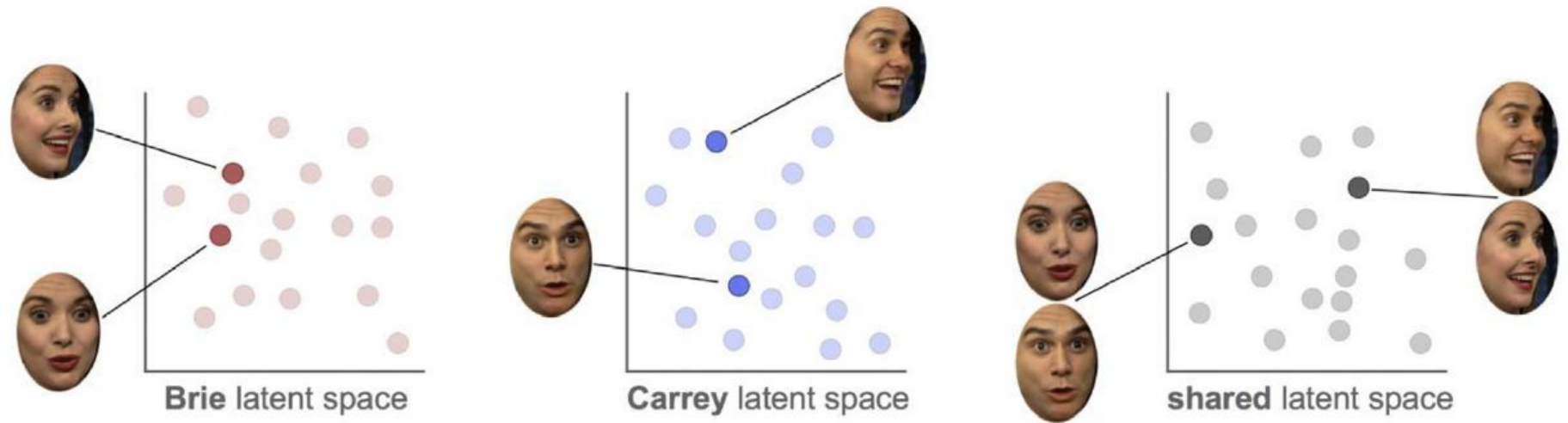




# How do deepfakes work?

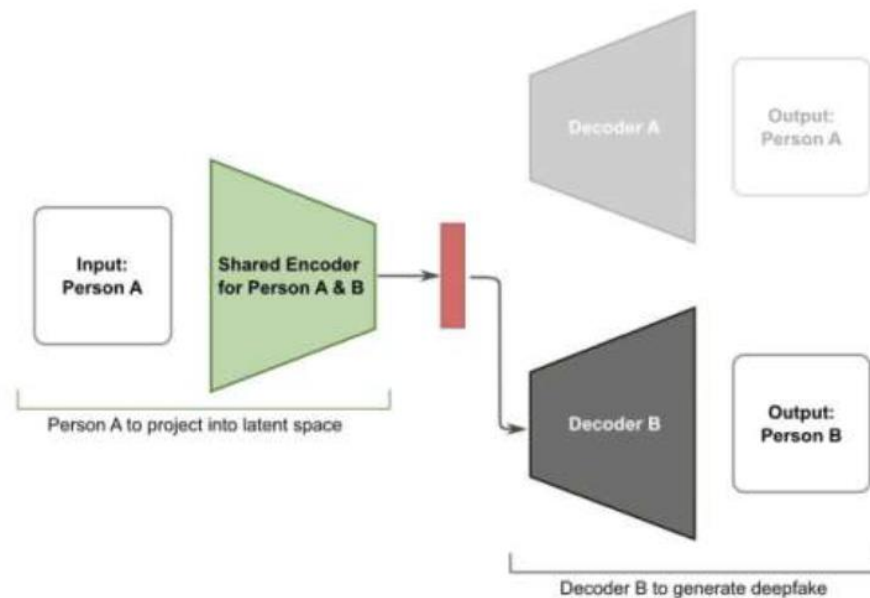
## The deepfake trick

Using the same encoder and hence latent space representation for images of two separate people



# How do deepfakes work?

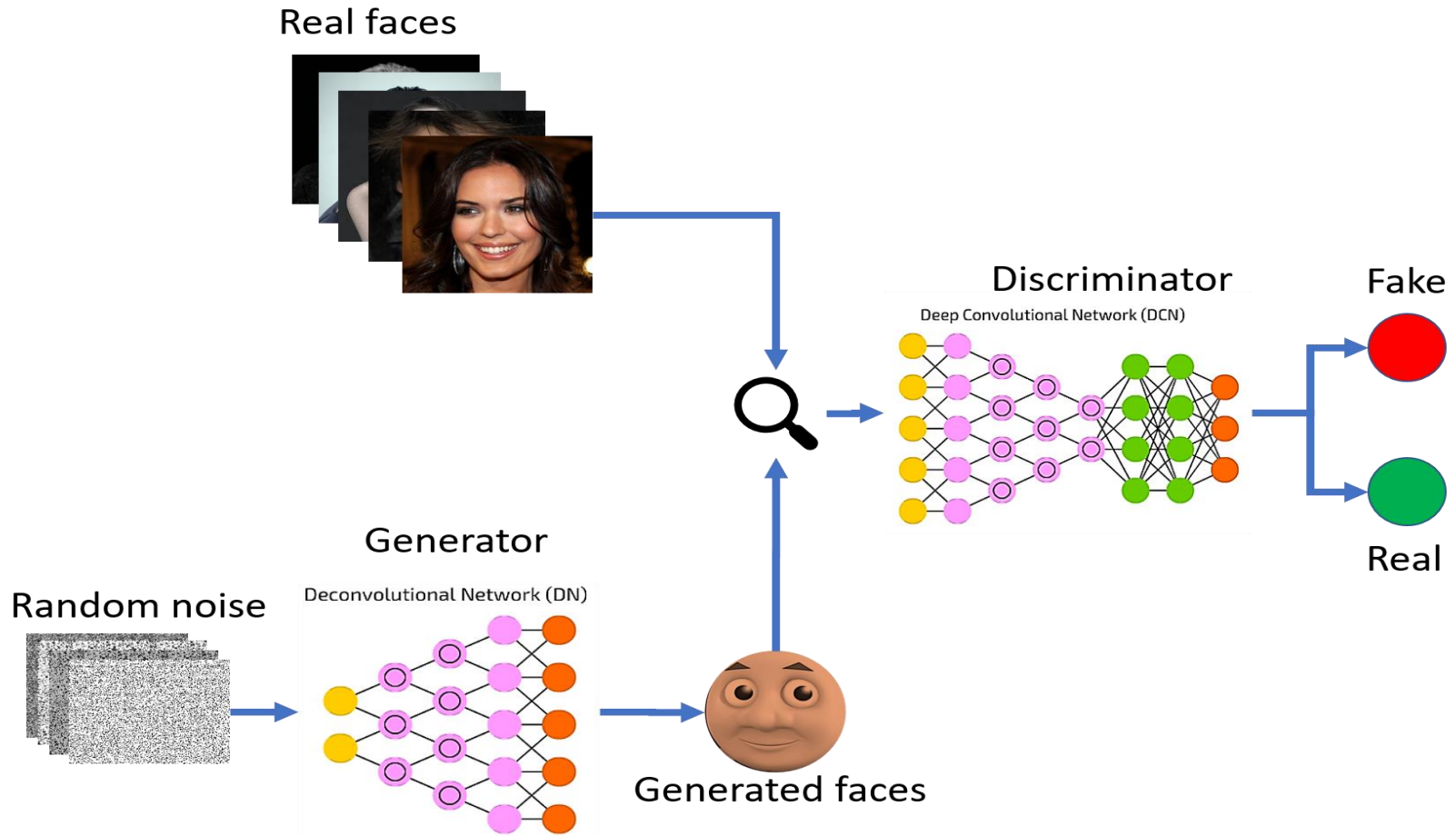
- A **shared encoder** is key to creating novel facial images of a target person that exhibit the same emotional expression, head posture, etc. as the original facial characteristics.
- This new image can then **be doctored back** into the original image to create a fake scene.



# Generative Adversarial Networks and Deepfakes

# Generative Adversarial Networks

[Goodfellow et al., 2014]



# Generative Adversarial Networks

## [Goodfellow et al., 2014]

---

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

---

**for** number of training iterations **do**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Sample minibatch of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{\text{data}}(x)$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$

**end for**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Update the generator by descending its stochastic gradient:

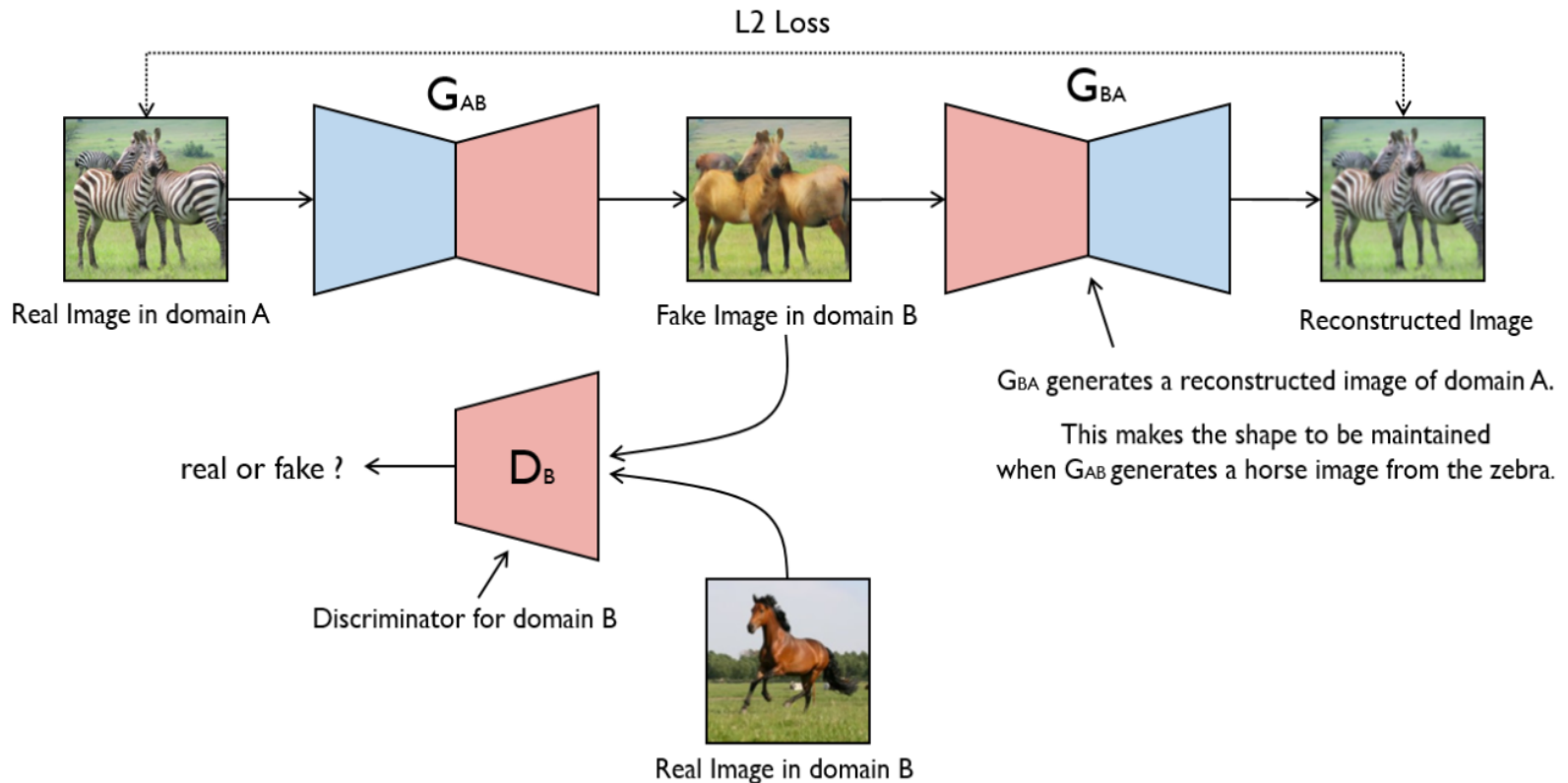
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

---

# CycleGAN



# GANs for good

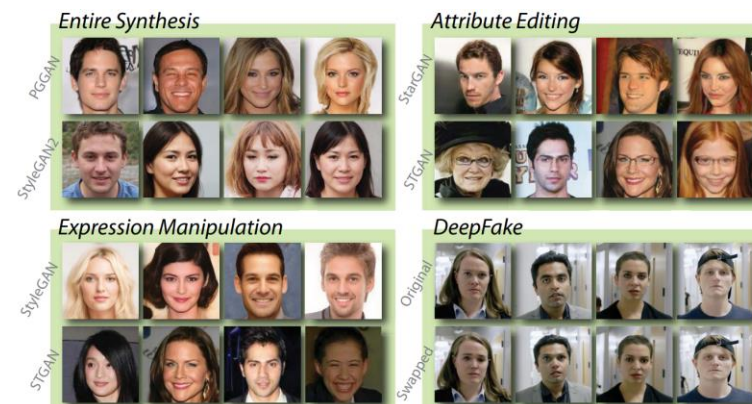
## Comprehensive generative objectives



Learning to Anonymize Faces for Privacy Preserving Action Detection, Ren et al. 2018

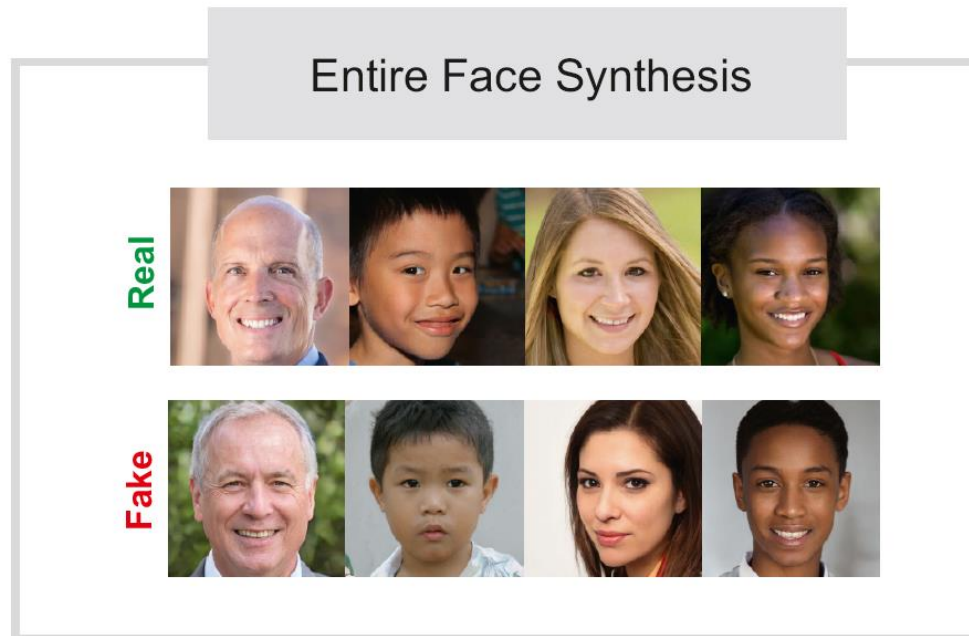
# GAN-based Attacks: four types

- i) Entire face synthesis
- ii) Identity swap (DeepFakes)
- iii) Attribute manipulation
- iv) Expression swap.





# Entire Face Synthesis



# Identity swap (DeepFakes)

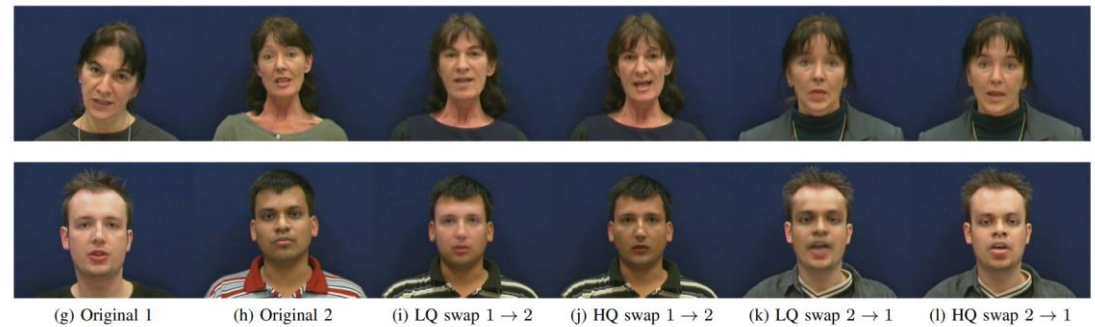
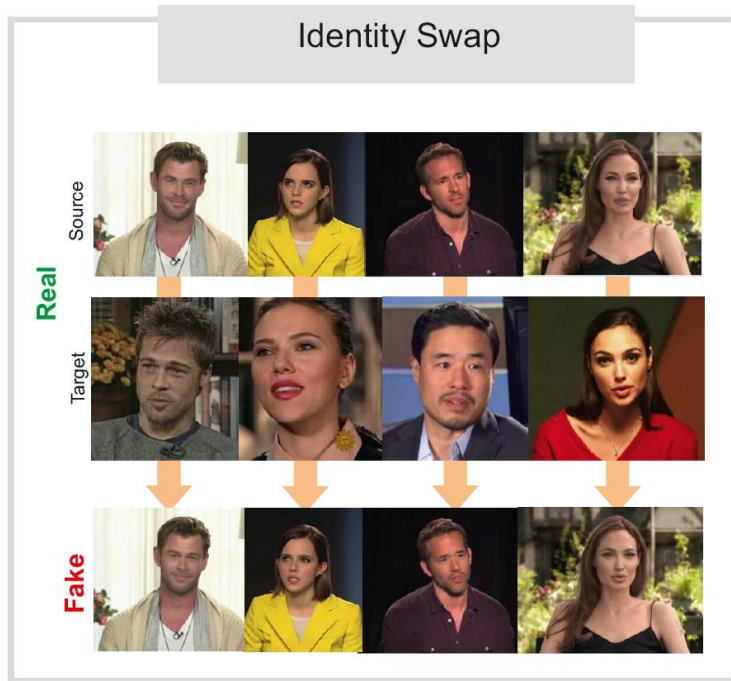
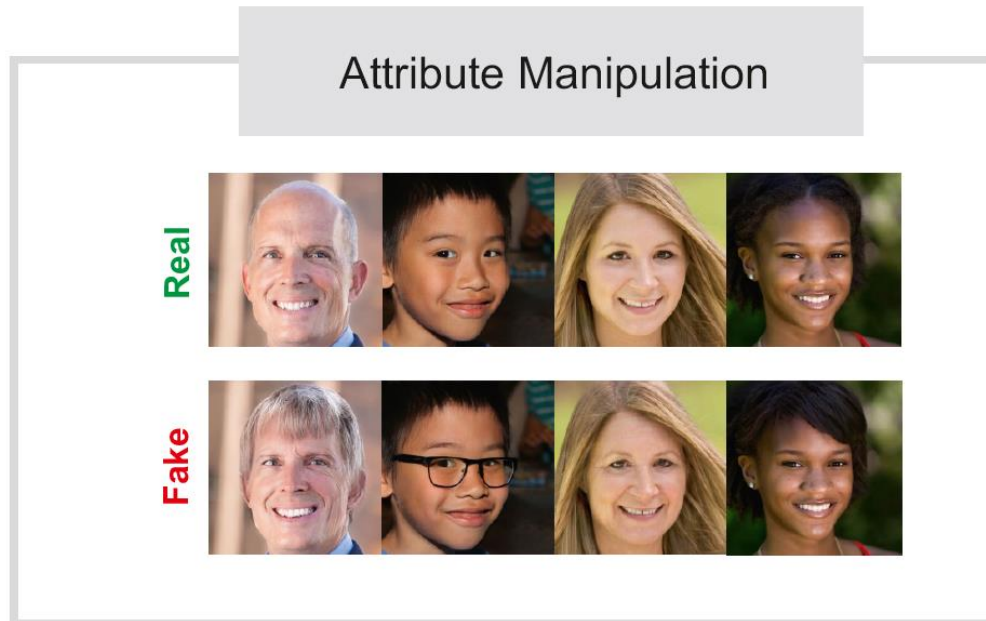
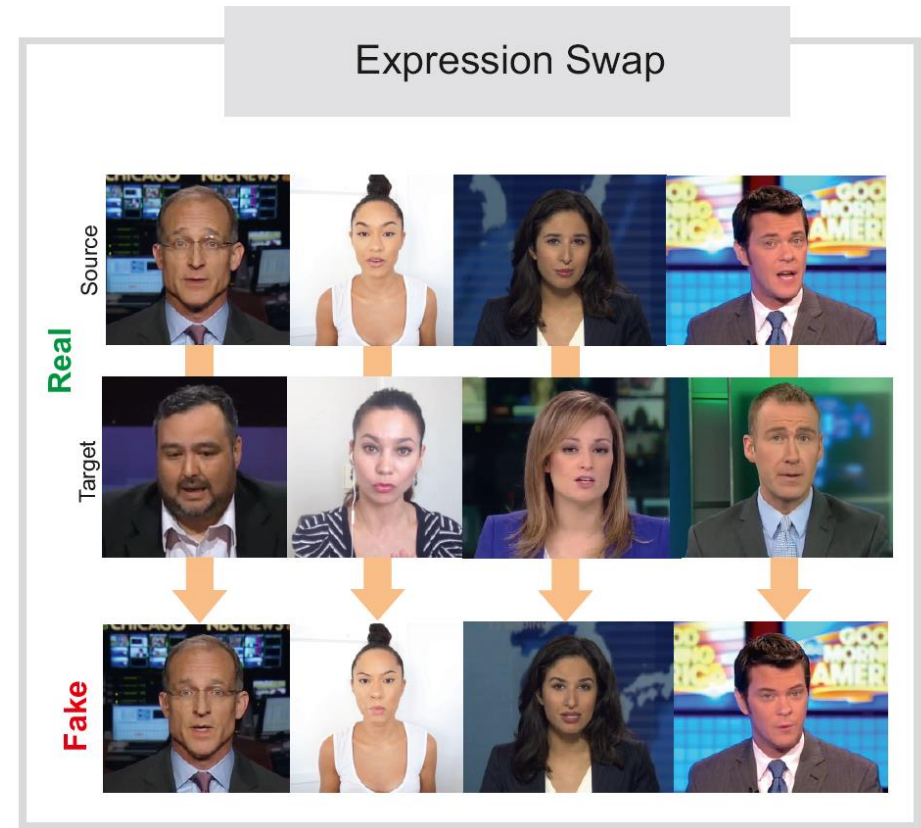


Fig. 1: Screenshot of the original videos from VidTIMIT database and low (LQ) and high quality (HQ) Deepfake videos.

# Attribute Manipulation

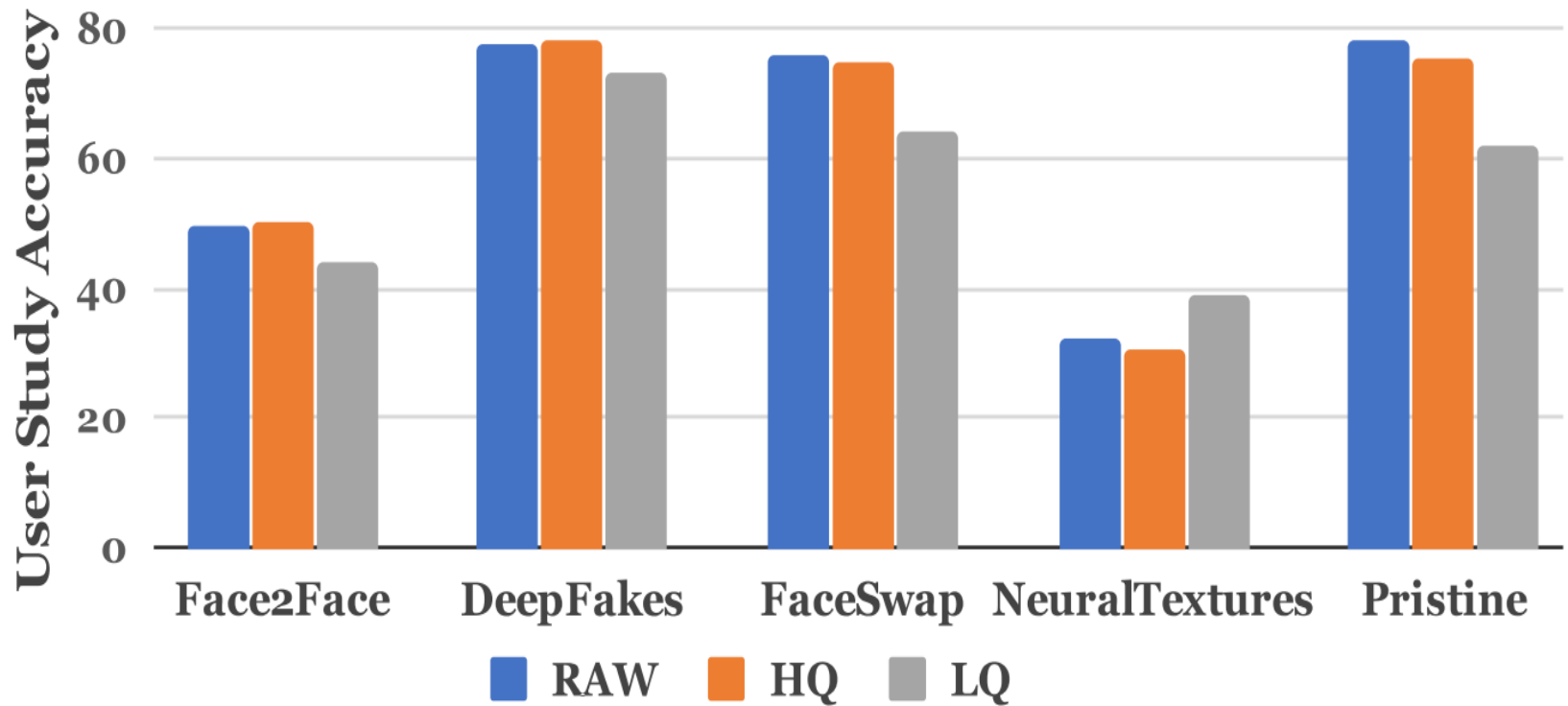


# Expression Swap



# How to detect Deepfakes

# Humans Can't Detect Image Manipulation Well



# Deepfake Detection Methodologies

- Signal level (sensor noise, double JPEG compression, etc.)
- Physical level (lighting conditions, shadows, reflections, etc.)
- Semantic level (consistency of meta-data)

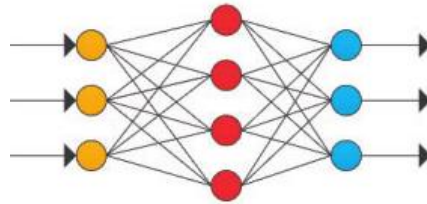
# Deepfake Detection Methodologies

- Physiological signals (breathing, pulse, eye blinking, etc.)
- Video authentication (ex. blockchain)



# Detection Approaches

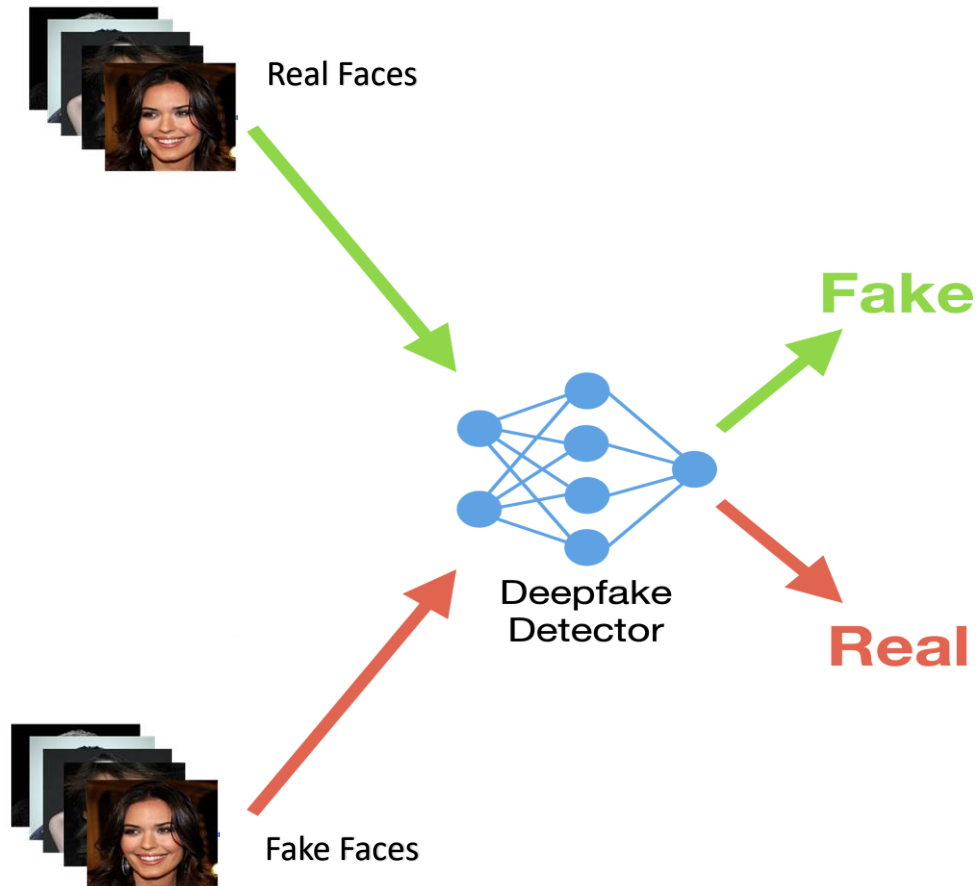
- Detection using Machine Learning/Deep Learning



- Detection using GAN fingerprints



# Machine Learning/Deep Learning



[cs.albany.edu/~lsw/celeb-deepfakeforensics.html](https://cs.albany.edu/~lsw/celeb-deepfakeforensics.html)

## Celeb-DF (v2): A New Dataset for DeepFake Forensics

Yuezun Li<sup>1</sup>, Xin Yang<sup>1</sup>, Pu Sun<sup>2</sup>, Honggang Qi<sup>2</sup> and Siwei Lyu<sup>1</sup>

<sup>1</sup> University at Albany, State University of New York, USA

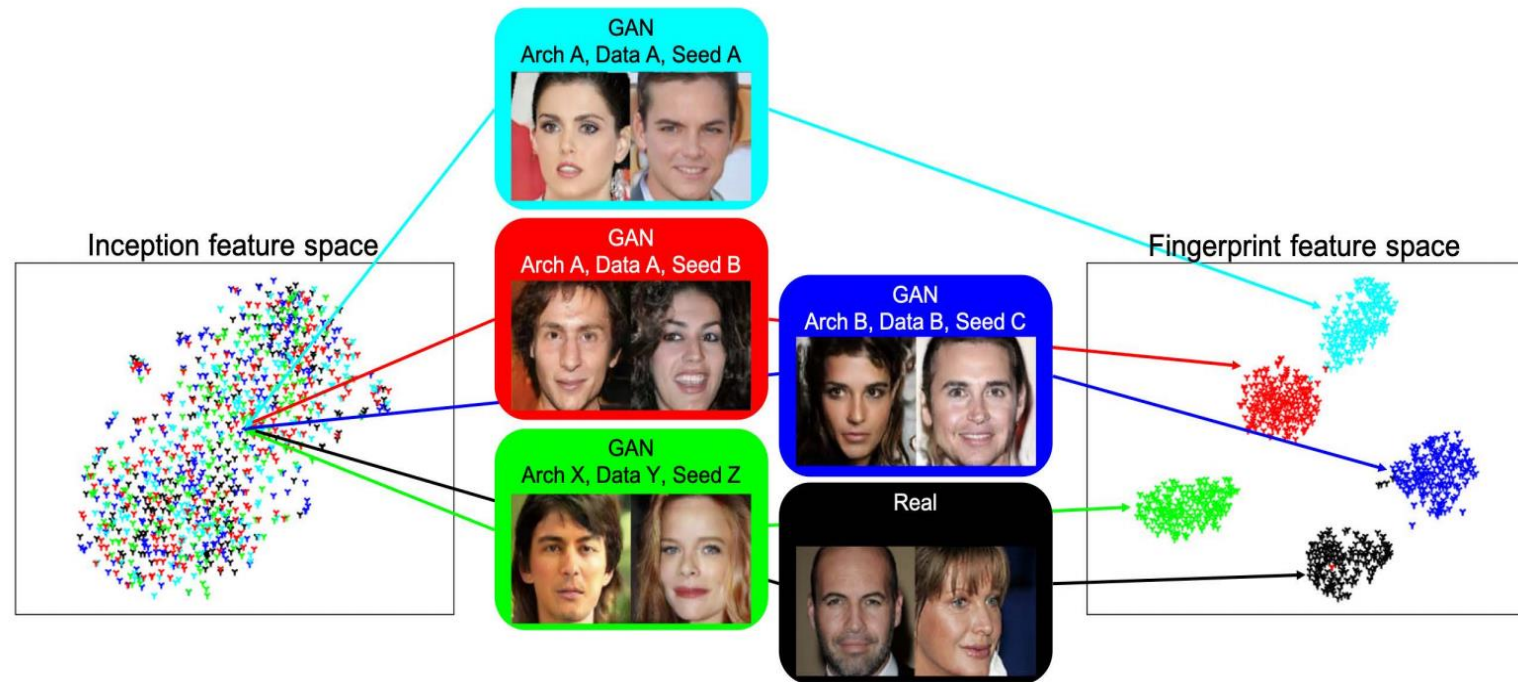
<sup>2</sup> University of Chinese Academy of Sciences, China

[Github](#) [Paper](#) [Celeb-DF \(v1\)](#)



Green box: Real images, Red box: Corresponding DeepFake images.

# Using GAN Fingerprint



# GAN-fingerprint removal approach

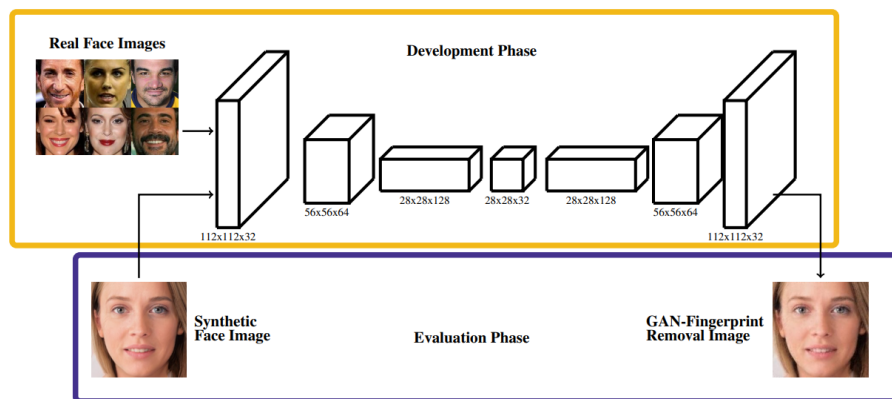
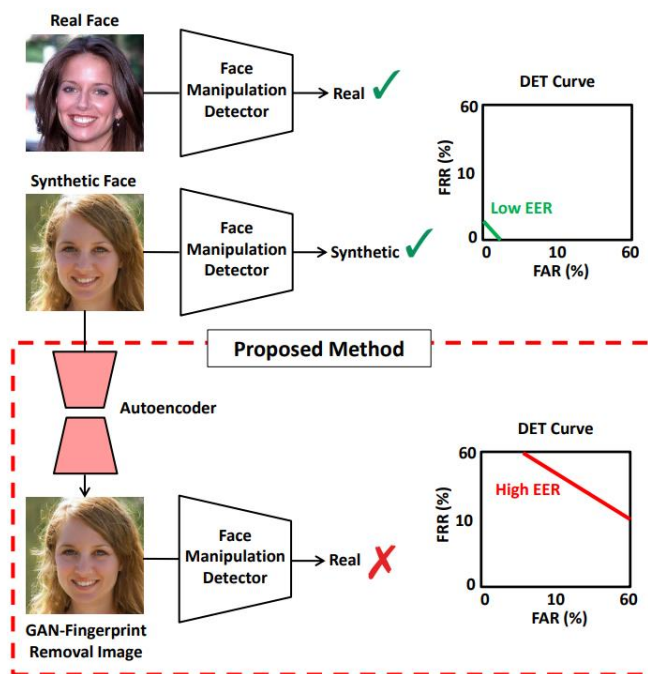


Fig. 2. **Proposed GAN-fingerprint removal approach based on convolutional autoencoders.** An autoencoder is trained with real face images from the development dataset. In the evaluation stage, once the autoencoder is trained, we can pass synthetic face images through it to remove the GAN fingerprint information, and also incorporate the correlations of real face images into the synthetic face images.



# Problem with current detection approaches

- most current face manipulations seem easy to be detected **under controlled scenarios**, i.e., when fake detectors are evaluated in the same conditions they are trained for.
- facial manipulation techniques are **continuously improving**.

These factors motivate further research on the generalization ability of the fake detectors against unseen conditions.

# The Real Business Impacts of Deepfakes

- False claims of malfeasance, damaging a product or company's reputation
- Endorsements that are not real (you thought fake written reviews were harmful)
- Video-backed HR complaints about a co-worker or a boss
- Insurance fraud, support by “video proof”
- False news about the company's owners, founders, leaders, etc.

# The Real Business Impacts of Deepfakes

- Onboarding processes subverted and fraudulent accounts created
- Identity theft, using video to convince someone to alter critical personal data
- Diversion of shipments
- Orders for unwanted materials
- Payments and/or funds transfer fraudulently authorized
- Blackmail based on the threat to release a damaging vid



# How to guard your organization against deepfakes

- **Employee training and awareness**
  - By offering adequate training and creating awareness employees can be turned into an additional line of defense.
  - Training should focus on how the technology is leveraged in malicious attempts and how this can be detected: enabling employees to spot deepfake-based social engineering attempts.

# How to guard your organization against deepfakes

- **Detection model**
  - Detecting false media **early** can help minimize the impact on your organization.
  - Developing new models that can **detect** fake images and videos.

# How to guard your organization against deepfakes

- **Response Strategy**

- Ensure that your organization is ready to adequately respond to a deepfake.
- Have a plan in place that can be set in motion when a deepfake is detected.
- It's important that individual responsibilities and required actions are defined in this plan.