

WROCŁAW UNIVERSITY OF SCIENCE AND TECHNOLOGY  
FACULTY OF ELECTRONICS

---

FIELD: Computer Science  
SPECIALIZATION: Internet Engineering (INE)

**MASTER OF SCIENCE THESIS**

Research on methods of changing objects in  
images using Deepfake technology

Badania metod zmiany obiektów na obrazach z  
wykorzystaniem technologii Deepfake

AUTHOR:  
Michał Zendran

SUPERVISOR:  
Dr inż. Andrzej Rusiecki

GRADE:

# Contents

<b>1</b>	<b>Abstract</b>	<b>2</b>
<b>2</b>	<b>Introduction</b>	<b>3</b>
2.1	Motivation . . . . .	3
2.2	Objective and assumptions . . . . .	4
2.3	State of the art . . . . .	4
2.4	Naming conventions and terminology . . . . .	4
<b>3</b>	<b>Theoretical background</b>	<b>5</b>
3.1	Artificial neural networks . . . . .	5
3.2	Convolutional neural networks . . . . .	6
3.3	Supervised, unsupervised and semi-supervised training . . . . .	6
3.4	Haar feature-based cascade classifiers . . . . .	6
<b>4</b>	<b>Deepfake methods</b>	<b>7</b>
4.1	Autoencoder . . . . .	7
4.2	Variational autoencoder . . . . .	8
4.3	VAE-GAN . . . . .	9
4.4	CycleGAN . . . . .	9
<b>5</b>	<b>Datasets</b>	<b>10</b>
5.1	Dataset description . . . . .	10
5.2	Data pre-processing . . . . .	11
<b>6</b>	<b>Technologies</b>	<b>12</b>
6.1	Software and Libraries . . . . .	12
6.2	Hardware . . . . .	12
<b>7</b>	<b>Networks implementation</b>	<b>13</b>
7.1	Variational auto encoder . . . . .	13
7.2	Convolutional variational auto encoder . . . . .	13
7.3	VAE-GAN . . . . .	13
7.4	CycleGAN . . . . .	13
<b>8</b>	<b>Results</b>	<b>14</b>
<b>9</b>	<b>Conclusions</b>	<b>15</b>
	<b>Bibliography</b>	<b>16</b>

# Chapter 1

## Abstract

To be written at the end of the work.

# Chapter 2

## Introduction

### 2.1 Motivation

Machine learning has found many, different applications in the field of image data processing and computer vision. From picture classification to image denoising and resolution enhancement, artificial neural networks has gained the opinion of exceptionally useful tool. But for some time, a new, controversial use-case has been getting more and more attention in both media and research circles. So-called "deepfake" technology has opened doors to many new possibilities of picture generation, but also raised many issues of moral and legal matters.

Deepfake is a technology from the field of machine learning designed to combine and overlay objects in images or videos creating deceptively realistic counterfeits. The name comes from combination of two terms: "deep learning" and "fake", and has its origins in a Reddit user named "deepfakes". Initially the term was associated only with face-swapping technology, but with time it was extended to all deep learning implementations of changing objects in images.

Deepfake technology has already found multiple applications such as changing seasons in the landscapes images, transforming horses into zebras or "repainting" images in styles of different artists. But the most controversial and impactful use-case so far is already mentioned face-swapping. In times of overwhelming amount of news it's getting harder and harder to filter out fake ones from valuable peaces of information. People generally tend not to check sources of information but rather blindly follow hot stories in social medias and television. Such environment combined with capabilities of deepfake gives possibilities of influencing elections by misrepresenting politicians in forged videos to defame or blackmail theme. Another popular use-case of described technology is creating erotic videos by replacing faces of porn actress with faces of well-known celebrities. This application might have less dangerous consequences than influencing world politics but may be hurtful to people that became objects of such act.

Although there are many malicious ways of using deepfake technology it might also be used for good reasons such as helping people to cope with the loss of the loved once or in entertainment filed by de-ageing actors to play younger-selves. Besides, to be able to detect and fight harmful applications of deepfake it might be vital to deeply understand algorithms and techniques behind it. Therefore, conducting research on that part of machine learning field seems to have great meaning in incoming times.

## 2.2 Objective and assumptions

This project aims to implement and compare four different methods of changing objects in images with application of artificial neural networks. For sake of this research, human faces were chosen as an object of replacement, as it rises a complex issue of simultaneous color, texture and shape modification.

As there are no numerical methods of measuring the quality of images obtained from deepfake algorithms, the only way of appraising results of methods discussed in this research is visual evaluation. To be able to fairly rate each implemented technique, the same set of images will be used as a learning dataset for all cases. Therefor, effects of all approaches will be visually evaluated and compared with each other, which will result in the final assessment. This rating of methods is the expected outcome of the research.

There are two main factors that will be taken into a consideration during a results evaluation process. First of them is a resemblance of the faked image, to the appearance of the imitated person. The more striking similarity, the better. The other crucial aspect is preservation of original facial expression and pose, as the believable deepfake must capture the source material movements. Resultant of those two factors will be the main feature to be rated. It is assumed that neither of mentioned characteristics should outweigh the other one, but rather, the final effect should be well-balanced composition of both aspects.

## 2.3 State of the art

While there are many, great articles and papers that elaborately explain all mentioned approaches of generating deepfakes, no comparison of those methods were found. To be written ...

## 2.4 Naming conventions and terminology

Below, all abbreviations and naming conventions used in this research are listed and explained:

- Deepfake – name of the deep learning technology of swapping objects in images or an end result generated with such technology
- ANN – artificial neural network
- CNN – convolutional neural network
- AE – autoencoder
- VAE – variational autoencoder
- GAN – generative adversarial network
- VAE-GAN – variational autoencoder-generative adversarial network
- Activation function – transfer function

# Chapter 3

## Theoretical background

### 3.1 Artificial neural networks

An artificial neural network is a computing system inspired by the neural structure of the biological brain and its way of processing information. Such structure is able to “learn” how to solve certain problems or recognize patterns without being pre-programmed with rules how to do it. Artificial neural networks are main tool used in deep learning, which is part of bigger family of machine learning.

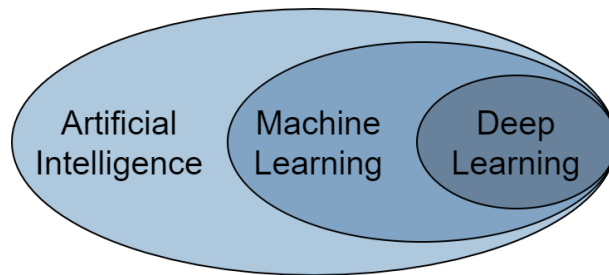


Figure 3.1: Deep learning belongingness

Fundamental unit of artificial neural network is so-called “artificial neuron” which is a mathematical function modeled after the structure of biological neuron. Such artificial neuron has a number of inputs  $\{x_0, x_1, x_2, \dots, x_n\}$  with assigned, separated weights  $\{w_0, w_1, w_2, \dots, w_n\}$  for each of them to be multiplied by. In the simplest case, all product are summed and passed to a transfer function  $\varphi$ . Additionally a bias value  $b$  is added to the sum of products which allows to shift the activation function up or down. Then the output  $y$  of the neuron is calculated as follows:

$$y = \varphi \left( \sum_{j=0}^n x_j w_j \right) \quad (3.1)$$

General idea of artificial neuron model is illustrated in figure 3.2. Typically, neurons are organized into layers which may perform different operations. The output of neurons from one layer is passed to the input of the next layer or might be a part of output vector of whole network. Modern models of artificial neural networks may consists of tens or even hundreds of different layers, on which depends how well given network will cope with certain tasks. Those layers, during process of network training, can adjust theirs weights to produce output values that correctly solve given problem.

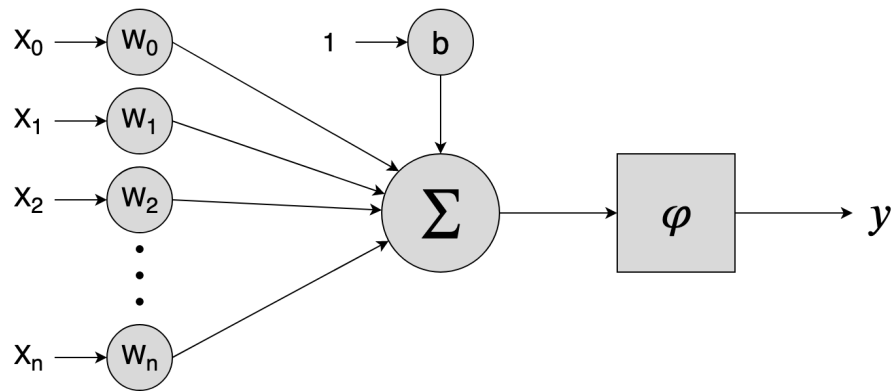


Figure 3.2: Artificial neuron model

## 3.2 Convolutional neural networks

Explain how it works, what are main use-cases and so on.

## 3.3 Supervised, unsupervised and semi-supervised training

Description of both and what are main differences and when to use which.

## 3.4 Haar feature-based cascade classifiers

As in section name

# Chapter 4

## Deepfake methods

### 4.1 Autoencoder

Autoencoders are the most basic approach to the problem of deepfake generation. In fact, all mentioned methods are just different variations of this idea. Autoencoder is a type of artificial neural network that learns to reproduce given input in an unsupervised manner. The problem is to train functions  $A : \mathbb{R}^n \rightarrow \mathbb{R}^p$  (encoder) and  $B : \mathbb{R}^p \rightarrow \mathbb{R}^n$  (decoder) to satisfy condition given in equation 4.1 as described in [1],

$$\arg \min_{A,B} E[\Delta(x, B \circ A(x))] \quad (4.1)$$

where  $E$ – expectation over the distribution of  $x$  and  $\Delta$ – reconstruction loss function, which measures the distance between given input and the output of the decoder. General idea of autoencoder model is illustrated in figure 4.1. Typically, architecture of autoencoder consists not only of input and output layers, as this would result in simple coping pixels from the input to the output of the network, but also contains single or multiple hidden layers in between, with the number of neurons lesser than the number of pixels in the input image. Such structure causes bottleneck effect and creates so-called compressed representation at the output of the encoder part, known also as “feature map” or in case of deepfake “latent face”. Such compression causes feature map to preserve only information most relevant for later reconstruction and gets rid of unnecessary data.

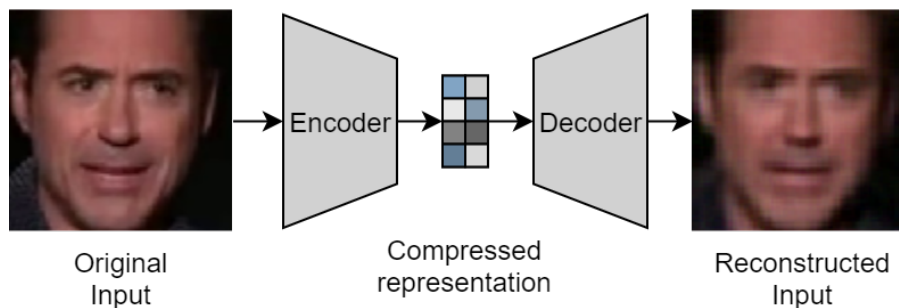


Figure 4.1: Autoencoder general idea

Generating deepfakes using autoencoders approach consists of three major steps. Described process is illustrated in figure 4.2. Let us assume that  $X$  is a set of face images of person  $x$  and  $Y$  is a set of face images of person  $y$ . Firstly, as shown in figure 4.2a, an encoder is trained to produce feature maps for images of face from both classes. Afterwards, two decoders are trained separately to reproduce original images from latent



faces generated by pre-trained encoder, as presented in figure 4.2b. Finally, decoders are switched to produce images from one class based on feature maps from the other class, which was illustrated in figure 4.2c.

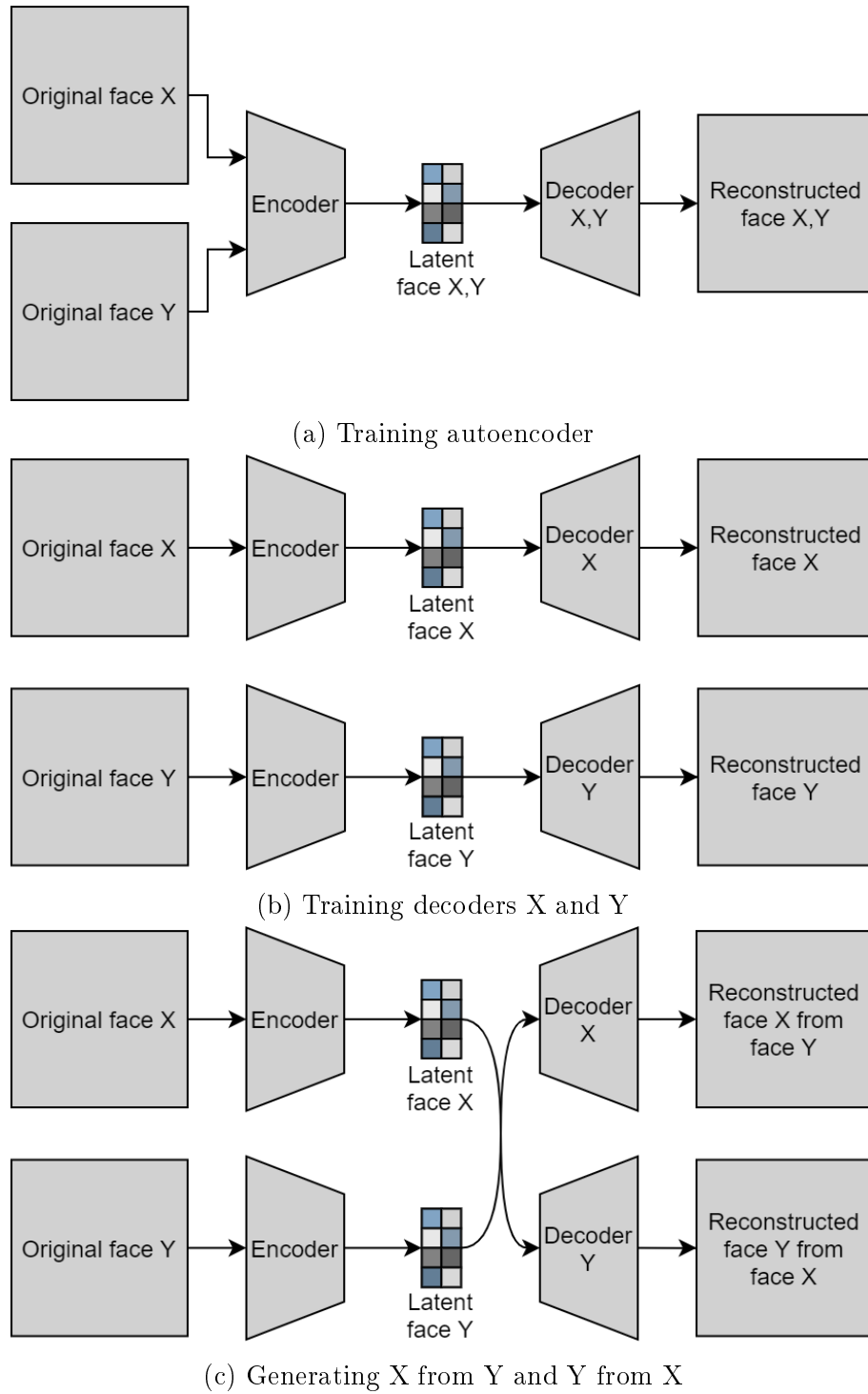


Figure 4.2: Three steps of generating deepfakes

## 4.2 Variational autoencoder

Idea behind deepfake generated by VAE with CNN

### 4.3 VAE-GAN

Idea behind deepfake generated by GAN actually "VAE-GAN".

### 4.4 CycleGAN

Describe what is it, what it consists of, what are its applications, why I thought it should work for deepfake. Explain how it works exactly. Show learning process and results (good ones: horses to zebras and bad ones: face to face). Idea behind deepfake generated by CycleGAN. Explain why I'm assuming it should it work?

# Chapter 5

## Datasets

### 5.1 Dataset description

Creating a high quality, well-balanced dataset is an essential step in the process of deepfake generation. Poorly prepared training set might cause that even great network architectures and state of the art algorithms will produce disappointing results. The most efficient approach to this problem ,in case of face-swapping technology, is obtaining images of targeted people from video recordings, which allows to produce great number of images that cover different facial expressions and head positions. There are several factors that should be taken into a consideration in order to construct proper dataset. First of all, videos of people with at least slight resemblance should be chosen. The most important feature in this matter is skin tone but the more similar appearances, the more deceiving deepfake might be achieved. Another important factor is overall quality of source material. Videos with high resolution and large number of frames per second are best suited for this purpose, but also require more computational power and time to properly train necessary networks. Finally, after a deconstruction of videos into single images, all set should be revised to remove pictures that depict faces which are blurry, deformed or in some ways covered, for example by hand.

For sake of this research the “VoxCeleb2” dataset was used. As described in [2], VoxCeleb2 consists of over 1 million utterances of over 6000 celebrities derived from videos from YouTube platform. Source data is diversified in terms of recorded people’s genders, ethnicities, ages and accents, but also in terms of videos quality, lighting conditions, stability and lengths. Each recording has a resolution of 224 by 224 pixels and depicts closeup shot of a character’s face and shoulders, which is close enough to clearly capture subtle facial expressions.



Figure 5.1: Sample images from VoxCeleb2 dataset

## 5.2 Data pre-processing

To prepare dataset best suited for purpose of this research following preparations were made, in accordance with guidelines described in section 5.1. Sample images of pre-processed data are presented in figure 5.2.

1. Two actors (Leonardo DiCaprio and Robert Downey Jr.), further called subject A and subject B, were chosen as targets of face replacement. This choice was dictated by how well their faces are known and recognizable which facilitates the final assessment and by factors mentioned in section 5.1.
2. From the set of all videos of chosen subjects, available in “VoxCeleb2” dataset, those which presented subjects in similar age and had good recording quality were selected.
3. From each video every tenth frame was extracted to limit the amount of nearly identical images. Additionally, through Haar feature-based cascade classifiers described in section 3.4, the face itself was cut out from each extracted frame to discard unnecessary parts of pictures.
4. Resolutions of all obtained face images varied, therefore data had to be rescaled to a common resolution of 160 by 160 pixels.
5. Final step, was to save obtained images as arrays into a single file in uncompressed “.npz” format. This allows easy data transferring and simplifies the process of loading data during the training of artificial neural networks.



(a) Sample images of subject A



(b) Sample images of subject B

Figure 5.2: Sample images of pre-processed data

# Chapter 6

## Technologies

### 6.1 Software and Libraries

As in title...

### 6.2 Hardware

As in title ... (My hardware, Google colab, Google cloud?)

# Chapter 7

## Networks implementation

Detailed description of implementation of each method. What are the topologies, what callbacks were used, why those parameter, why those batches itp

### 7.1 Variational auto encoder

### 7.2 Convolutional variational auto encoder

### 7.3 VAE-GAN

### 7.4 CycleGAN

# Chapter 8

## Results

Presentation and discussion of results for each method

## Chapter 9

## Conclusions



# Bibliography

- [1] Dor Bank, Noam Koenigstein, and Raja Giryes. Autoencoders, 2020.
- [2] J. S. Chung, A. Nagrani, and A. Zisserman. Voxceleb2: Deep speaker recognition. In *INTERSPEECH*, 2018.

# List of Figures

3.1	Deep learning belongingness . . . . .	5
3.2	Artificial neuron model . . . . .	6
4.1	Autoencoder general idea . . . . .	7
4.2	Three steps of generating deepfakes . . . . .	8
5.1	Sample images from VoxCeleb2 dataset . . . . .	10
5.2	Sample images of pre-processed data . . . . .	11

# List of Tables