
The generalization challenge in deepfake image detection on human faces using generative adversarial networks (GANs)

1
2
3
4

Anonymous Author(s)

5
6
7
8
9
10
11
12
13
14
15
16

Abstract

Over the past few years, alongside the rising popularity of AI technology, deepfakes did become prevalent in our lives. Images that are generated by Generative adversarial networks can now easily fool a human due to their high quality. This comes with multiple issues, when the technologies are misused for the creation of fake news or false identities. These are problems that no one quite seemed to be able to push back. Many researchers all over the world have approached this topic but failed to build applicable solutions for our everyday lives, while deep fakes on the other hand are still improving. We propose a project to support the research in developing a well generalizing GAN that will be able to not only produce and classify training images of human faces, but also adapt to various settings using large scale datasets.

17

1 Introduction.

1.1 Problem definition

As technologies become more advanced, humans more often lose the competition against the machines. In some fields, we have now reached a point where machines can easily do tasks humans cannot. Looking at the two pictures (Figure 1 and 2), the question can be asked, whether we can distinguish which image is AI-generated and which one is real.

25



26
27

Figure 1: Human face¹



Figure 2: Human face

¹ Both images:

<https://www.nytimes.com/interactive/2024/01/19/technology/artificial-intelligence-image-generators-faces-quiz.html>

28 Most people will not guess the right solution. Both pictures above are generated by AI.
29 AI filters and deepfake have been developed to the point where we can no longer identify
30 them with our own eyes [1]. Many of these technologies are very popular on social media
31 these days because they are easily and widely available and make it simple to edit or generate
32 more popular posts [2]. Several factors are increasing the widespread of deepfakes. For
33 example, open source deepfake generation tools make it easier to use deep fake technology
34 through mobile phone apps by lowering barriers to creating manipulated media. In addition,
35 various online communities and companies have provided a platform to share technology
36 with deep fake media, fueling the spread of deepfake [1, 3].

37 However, many issues that exploit false voices, images, and deepfakes tend to be dangerous
38 in many areas of our lives. There are some examples that show that it can affect everybody.
39 In some cases, even US presidents have been victims of deep fakes reaching from funny
40 videos to serious political matters [4]. To name another celebrity that became a victim of
41 deepfakes, in March 2023, a video was released that was manipulated using deepfake where
42 Bill Gates ended the interview, embarrassed by a question about COVID vaccine during a
43 news show interview [5]. Furthermore, most of the current legal and ethical concerns
44 regarding synthetic media have focused on the issue of “deepfake pornography” or the
45 non-consensual distribution of synthetically generated images, which is usually associated
46 with an individual's face being mapped onto pornography content without the individual's
47 knowledge [6]. Additionally, individuals and criminal organizations have already started
48 utilizing the widely available systems to commit different crimes that involve identity theft,
49 such as financial fraud [2, 3, 6].

50 As such, deepfake is increasing several risks, such as political disinformation, identity and
51 financial fraud as well as deepfake pornography. Although most deepfake is still associated
52 with funny celebrity videos, it opens opportunities for dangerous new fraud and slander
53 because it evolves too quickly to be caught by authorities missing advanced protection
54 systems [1-13]. As Generating software is widely available there is a lack of protection
55 methods not only for individuals but also organizations with good financial resources,
56 showing there has not yet been found a way to effectively fight the issues arising from the
57 fast-going development in this field.

58

59

60 1.2 Literature review

61 During the past few years, deepfake generating systems have been developed many attempts
62 to improve the detection of synthetic images. Most popular nowadays are so-called GANs, or
63 Generative Adversarial Networks, that consist of a generator and a discriminator [2, 6-13].
64 During training processes, while the generation of synthetic images improves, the detection
65 of those between natural images also improves. The improvement of these systems has
66 become a hot topic not only in society but also in research.

67 Preeti et al. explained that not only manipulation of faces, but full identity swaps can
68 nowadays be done with open-source software, allowing the public to generate high-quality
69 and realistic synthetic images of human faces in a timespan of only a few seconds. They
70 explain that GANs, like the so-called StyleGAN, make this real-world application possible.
71 Furthermore, they argue that this rapid development and easy access to the public are
72 enabling criminal misuse, such as financial fraud. They continue by stating that the generated
73 images from modern GANs are so advanced that even humans cannot distinguish between
74 real and fake images, leaving room for opportunities to come up with ways to
75 computationally distinguish between them. They emphasize that among the many
76 well-working generating models, there are only a few that bring up a discriminator that is
77 able to perform well on new unseen data and that it is important to now build generalizing
78 solutions to fight deepfake-based crimes [2].

79 Almar's study presents diverse deep learning approaches for detecting GAN-generated
80 content, where deep learning has notably excelled. Techniques include CNN, RNN, long
81 short-term memory (LSTM), a type of artificial recurrent neural network, and methods for
82 analyzing statistical features of images using image preprocessing techniques and improving

83 the detection of fake face images. He explains that each of these implementations of GAN
84 works well on either images or videos during training and validation, while suffering from a
85 lack of good performance on new data. Additionally, he introduces Forensic CNN, which
86 employs Gaussian Blur and Gaussian Noise for better model generalization. However, this
87 might bring small improvements; generalization remains challenging. He states that with the
88 mass of new generated images every day, there is still an issue finding well-curated datasets,
89 which might be another reason for not much visible improvement regarding the
90 generalization of the introduced models and that this as well requires further development
91 [7].

92 A team from Catania recently mentioned that GANs leave identifiable traces determined by
93 their architecture and specific parameters when generating synthetic images. They examined
94 whether methods that utilize frequency domain analysis to detect these traces perform well
95 compared to augmentation methods. For instance, some detectors utilize the discrete cosine
96 transform (DCT) to analyze images or extract features from DCT blocks to identify unique
97 traces. In the challenge, the focus was on evaluating the detectors' resilience to alterations in
98 images, with analytical approaches based on DCT analysis demonstrating the highest
99 accuracy. After comparing the method to a synergy of CNN and Vision Transformers to allow
100 large-scale classifications with both local and global vision, they found out that the
101 application of DCT achieved high accuracy while employing an analytical approach, even
102 with complex datasets. An important point they mentioned aside from the generalization
103 issue, which was reduced by regularization, was that the analytical method also improved the
104 explainability of the classification result. They explained that explainability is another
105 important part, since humans can no longer identify fake images and are more likely to trust
106 a system they understand [8].

107 In comparison to previous ones, Zha et al. introduced a different approach to solve the issues
108 that appeared alongside the development of generative models. They explained that due to
109 the nature of most internet platforms, videos and images get compressed while being
110 uploaded, which leads to fading or loss of image attributes that are often the ones GAN
111 discriminators focus on for their detection, leading to them not being useful in a real-life
112 setting. The idea they proposed is to implement a model that focuses on the robust
113 representations of each class using a real-centric hard feature fusion method (footnote that
114 says: important for pattern classification using fine-grain features) to distinguish inter- and
115 intra-class features. Even though they were not able to achieve results using this rather
116 experimental approach further away from the state-of-the-art method, studies like this show
117 there is still demand for new and different ways of tackling the problem. Even though they
118 proposed a different method, they still emphasized the importance of improvement of
119 generalization throughout the field [9].

120 As we can see throughout the field, most researchers agree that the main problem that needs
121 to be solved is the lack of generalizing models to build useful systems to reliably detect
122 deepfakes [1-3, 6-13].

123 1.3 Project objectives

125 Combining the earlier statements, our goal is to build a well-generalizing classifier for
126 synthetic and real images. As we have seen, it is broadly agreed upon that with the
127 development of models that generate images so close to real images that humans cannot
128 detect the differences anymore, there is a demand for lasting solutions for systems that are
129 able to consistently detect these deepfakes in order to prevent negative consequences of the
130 achieved development [1-11]. There is, therefore, a need for solutions that outperform
131 humans more and more. Our goal is, therefore, to build a system that not only is able to
132 generate excellent fake images but also detect these, while humans cannot. Additionally, it
133 will not be enough to build another GAN that generates images well. We need to find a way
134 to train the discriminator inside the GAN so that it will not only detect images generated by
135 the corresponding generator but also those created through different systems and of different
136 quality.

137

138 **2 Scope of work**

139 The scope of work in this project will be limited to the discussed issue of creating well
140 generalizing solutions. As the title already states, the project includes the task to set up a
141 convolutional GAN model. The GAN architecture includes a Generator that generates new
142 synthetic images from a random input and typically gets better during training while it tries
143 to fool the second part of a GAN, the discrimination model. To build the model, whose
144 performance on new data you want to later improve the discrimination between
145 well-generated synthetic images and the natural training data, the discriminator should be a
146 convolutional architecture that works well on classifying images, such as ResNet or other
147 advanced solutions [10]. Due to limited time and the availability of many tutorials and
148 pretrained models, it is suggested to use pretrained models and then working on top, making
149 adjustments to the architecture and trying techniques to improve generalization. However, if
150 someone wants to implement it on their own, that is acceptable as well. Whatever way the
151 model is implemented, whether it is pretrained or build and trained from scratch, it needs to
152 be runnable in Colab.

153 To train the model on classifying synthetic and real images, the synthetic images should be
154 generated by the generator. Therefore, for the training only natural non synthetic images need
155 to be used to make sure the model is not trained to identify generated images as natural faces.
156 There is a dataset of natural faces collected by Nvidia that is accessible on Kaggle via
157 <https://www.kaggle.com/datasets/xhlulu/140k-real-and-fake-faces?select=test.csv> . This
158 dataset should be used as your training dataset. In case of using additional data, please refer
159 from using the one provided on <https://www.chicagofaces.org/download/> , since it will be
160 used for the evaluation of the generalization and therefore should not be used for training
161 purposes.

162 **3 Evaluation criteria**

164 The main evaluation criteria for this project will be the accuracy on the provided test dataset.
165 Since we are trying to tackle the issue of the lack of generalization in current methods, the
166 performance on new, unseen data is what we are interested in when evaluating our results.
167 This is further the case, because deepfake generation itself is advancing at a fast pace, and a
168 model that does not generalize well on existing unseen data will most likely not be able to
169 detect even more advanced fakes, therefore having no chance of making an impact in the
170 field of deepfake detection [2,11-13]. Furthermore, computational cost and test runtime are
171 important to manage, not only because we are working with limited Colab resources and
172 limited time, but also because if we want to apply our solution to a social media platform
173 with millions of new images uploaded every day, the system needs to take as little time and
174 computation as possible to be useful in the real world [2,9]. For many other issues like fraud
175 detection, fast and cost-effective computation is also mandatory to catch the criminals before
176 damage has been done to the victim. In conclusion, the evaluation should be done on the test
177 accuracy with respect to runtime on the test dataset.

179 **References**

- 180 [1] Dr.A.Shaji George, A.S.Hovan George, “Deepfakes: The Evolution of Hyper Realistic
181 Media Manipulation”, Partners Universal Innovative Research Publication (PUIRP), vol 01,
182 PU Publications, 2023, pp 58–74
- 183 [2] Preeti et al., “A GAN-Based Model of Deepfake Detection in Social Media”, Procedia
184 Computer Science, 2023
- 185 [3] Bateman, Jon, “Deepfakes and synthetic media in the financial system: Assessing threat
186 scenarios”, Carnegie Endowment for International Peace, 2022
- 187 [4] Helmus, Todd C., “Artificial Intelligence, Deepfakes, and Disinformation: A Primer”,
188 *Jstor*, 2022.
- 189 [5] Jack Goodman, “Coronavirus: The Fake Bill Gates Post and Other Claims to Ignore”,
190 *BBC News*, BBC, 2020, www.bbc.com/news/52039642.
- 191 [6] Raza, Ali et al., “A Novel Deep Learning Approach for Deepfake Image Detection”,
192 Applied Sciences, 2022.
- 193 [7] Almars, Abdulqader M, “Deepfakes Detection Techniques Using Deep Learning: A
194 Survey”, *Journal of Computer and Communications*, 2021
- 195 [8] Guarnera, Luca et al., “The Face Deepfake Detection Challenge”, *Journal of Imaging* 8,
196 2022
- 197 [9] Zha, Ruiqi et al., “Real-centric Consistency Learning for Deepfake Detection”, *ArXiv*,
198 2022
- 199 [10] Guarnera, Luca et al., “Level Up the Deepfake Detection: a Method to Effectively
200 Discriminate Images Generated by GAN Architectures and Diffusion Model.” ArXiv, 2023
- 201 [11] Nguyen, Thanh Thi et al., “Deep learning for deepfakes creation and detection: A
202 survey”, *Comput. Vis. Image Underst.* Vol. 223, 2022
- 203 [12] Chen, Liang et al., “Self-supervised Learning of Adversarial Example: Towards Good
204 Generalizations for Deepfake Detection”, IEEE/CVF Conference on Computer Vision and
205 Pattern Recognition (CVPR), 2022
- 206 [13] Xiao, Shuai et al., “MCS-GAN: A Different Understanding for Generalization of Deep
207 Forgery Detection”, IEEE Transactions on Multimedia 26, 2024

208 **Training Dataset**

209 140k Real and Fake Faces [Data set,] Kaggle, 2020

210 **Test Dataset**

211 Synthetic Faces High Quality (SFHQ) part 1 [Data set], Kaggle, 2022

212 Chicago Face Database, CFD, 2023