

Exploring the Adversarial Robustness of AI-generated Image Detectors

Thomas Lazzerini, Samuele Cappelletti, Martina D'Angelo
University of Trento

Abstract—TODO

I. INTRODUCTION

Synthetic images are now flooding the real world. From online dating sites to social media, fake profiles and scams are everywhere. The problem with synthetic images is that, while some of them are funny and harmless, others could be harmful, they could be exploited by malicious users [1]. In relation to this, in the image forensic field there is a continuous fight between *fake image detectors* and *adversarial attacks*. On one hand, the detectors try to distinguish fake images from real ones, while, on the other hand, the attacks try to trick the detectors by manipulating the images (both real and fake ones). In order to detect fake images, we can exploit the traces/artifacts that fake image generators leave on the generated images. To do so, we have two main types of techniques: the *low-level forensic techniques* and the *high-level forensic techniques*. The former focuses on the pixel-level artifacts, which are almost invisible to the human eye. The latter focuses on physical inconsistencies and on repeated and uniform patterns, both of which are mostly visible to the human eye. Examples of physical inconsistencies are lighting, shadows, reflections or vanishing points inconsistencies [2][3]. While, an example of repeated and uniform patterns, typical of GAN-based image generators, is the generation of the mouth, the nose and the eyes always in the same position [4]. In general, we prefer to rely on low-level artifacts since fake image generators are becoming smarter every day, thus they are learning to generate always more realistic images, with fewer physical inconsistencies.

II. DETECTORS

In this section we will briefly describe a couple of fake image detectors: one uses CLIP to extract the feature vectors from the images [5] and one identifies the low-level traces/artifacts by training a GAN and a Diffusion Model [6].

A. CLIP-Based Detector

Many SoTA fake image detectors work very well in detecting fake images that are generated by an image generator of the same family of the generator that generated the images that they were trained on. But, the problem is that their performance decreases a lot when trying to detect images generated by another type of detector. For example, if the images used to train the detector were generated by a GAN-based generator, then the detector is good in detecting images

generated by other GAN-based generators but, it is bad in detecting images generated by a Diffusion-based generator. Moreover, the SoTA detectors hardly generalize to new and unseen generative methods.

On the other hand, the CLIP-based detector proposed by [5] works well in detecting images generated by any type of generator, both with and without augmentations (e.g., cropping, resizing, compression, etc.). Moreover, the performance of this CLIP-based method is similar to the one of the SoTA detectors in the in-distribution scenario but, it has a significant improvement in the out-of-distribution scenario. This is thanks to the fact that the CLIP features achieve an excellent generalization and robustness even with a few examples (e.g., 1k or 10k).

The CLIP-based method consists in: collect N real images $\{R_1, \dots, R_N\}$ with their corresponding captions $\{t_1, \dots, t_N\}$. Then, use these N captions to generate N images $\{F_1, \dots, F_N\}$ using some image generator $G(\cdot)$, $F_i = G(t_i)$. Successively, use CLIP to extract the feature vectors $\{r_1, \dots, r_N\}$ and $\{f_1, \dots, f_N\}$ of the $N + N$ images (real + generated), $r_i = \text{CLIP}(R_i)$ and $f_i = \text{CLIP}(F_i)$. Finally, feed the $N + N$ feature vector to a linear SVM classifier.

In the paper they compared the results (in AUC) of the CLIP-based method with other SoTA detectors on images generated by different generators of different families: *GAN*, *Diffusion* and *Commercial Tools*, with and without post-processing on the images. In the latter, the performance of the CLIP-based method is the best one in average wrt. the performance of the other SoTA methods tested. In the case of post-processed images the results are slightly worse but, especially for the CLIP-based method, they are still good across all generators.

B. PIZZA

III. ATTACKS

A. Mimicry

B. SD Laundering

C. White Black

D. Adversarial Robustness

IV. EXPERIMENT

V. CONCLUSIONS

REFERENCES

- [1] N. Carlini and H. Farid, "Evading deepfake-image detectors with white- and black-box attacks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 658–659.

- [2] H. Farid, “Lighting (in) consistency of paint by text,” *arXiv preprint arXiv:2207.13744*, 2022.
- [3] —, “Perspective (in) consistency of paint by text,” *arXiv preprint arXiv:2206.14617*, 2022.
- [4] S. Mundra, G. J. A. Porcile, S. Marvaniya, J. R. Verbus, and H. Farid, “Exposing gan-generated profile photos from compact embeddings,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 884–892.
- [5] D. Cozzolino, G. Poggi, R. Corvi, M. Nießner, and L. Verdoliva, “Raising the bar of ai-generated image detection with clip,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 4356–4366.
- [6] R. Corvi, D. Cozzolino, G. Zingarini, G. Poggi, K. Nagano, and L. Verdoliva, “On the detection of synthetic images generated by diffusion models,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.