

Exploring the Adversarial Robustness of CLIP for AI-generated Image Detection

Vincenzo De Rosa, Fabrizio Guillaro, Giovanni Poggi, Davide Cozzolino and Luisa Verdoliva
University Federico II of Naples, Italy

Email: {vincenzo.derosa3, fabrizio.guillaro, poggi, davide.cozzolino, verdoliv}@unina.it

Abstract—In recent years, many forensic detectors have been proposed to detect AI-generated images and prevent their use for malicious purposes. Convolutional neural networks (CNNs) have long been the dominant architecture in this field and have been the subject of intense study. However, recently proposed Transformer-based detectors have been shown to match or even outperform CNN-based detectors, especially in terms of generalization. In this paper, we study the adversarial robustness of AI-generated image detectors, focusing on Contrastive Language-Image Pretraining (CLIP)-based methods that rely on Visual Transformer backbones and comparing their performance with CNN-based methods. We study the robustness to different adversarial attacks under a variety of conditions and analyze both numerical results and frequency-domain patterns. CLIP-based detectors are found to be vulnerable to white-box attacks just like CNN-based detectors. However, attacks do not easily transfer between CNN-based and CLIP-based methods. This is also confirmed by the different distribution of the adversarial noise patterns in the frequency domain. Overall, this analysis provides new insights into the properties of forensic detectors that can help to develop more effective strategies.

Index Terms—Deepfakes, AI-generated image detection, adversarial robustness.

I. INTRODUCTION

AI-generated images are here to stay. People no longer pay much attention to the nature of an image, whether it is generated by a conventional device, created by a neural network or, more often, acquired by a smartphone with plenty of AI-based enhancement filters. However, there are many situations, in journalism, politics, or the judiciary, where establishing the nature of an image, real or synthetic, is of great importance. In recent years there has been intense research on this topic. A number of CNN-based detectors have been proposed which can easily recognize a synthetic image when it is generated by an AI model seen in the training phase. Unfortunately, their performance degrades sharply on images generated by new models (not an uncommon case) or impaired by compression or resizing [1]. Very recently, however, several new detectors have been proposed based on transformer architectures or that rely on features extracted from CLIP [2]–[5]. These approaches show very promising result in terms of generalization to unseen synthetic samples. This is in-line with current findings in computer vision that state that the transformer’s self-attention-like architecture is a key ingredient for improving the performance on out-of-distribution samples [6].

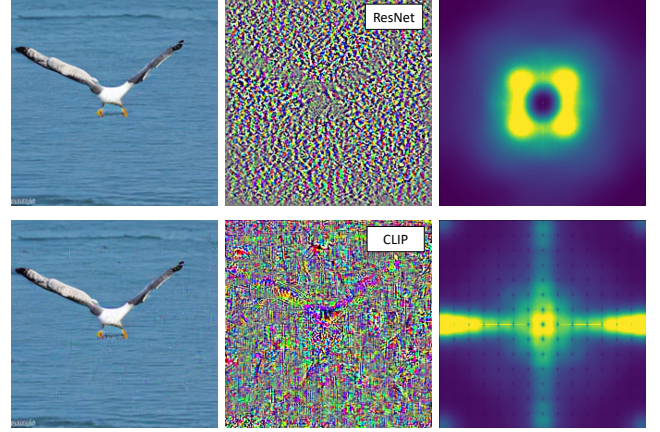


Fig. 1: **CLIP vs. ResNet.** l_2 -PGD attacks to ResNet-based (top) and CLIP-based (bottom) detectors. From left to right, attacked image, magnified adversarial perturbation, average spectrum of the adversarial noise. Attacks to Transformer-based detectors work on lower frequencies than attacks to CNN-based detectors do. In addition they present a clear cross-shaped directional spectrum, as also shown in [14] for image classification, which is due to the patch-wise processing before the transformer blocks.

In this work we want to investigate the robustness to adversarial attacks of such powerful forensic detectors and compare them with CNNs. Adversarial deep learning is a thriving field of research in its own right. This is even more true for forensic detectors [7]. In fact, to distinguish real images from those generated by AI they rely on subtle and easily perturbed traces [8], a sort of artificial fingerprints that characterize generative architectures [9]–[11]. Early works have demonstrated that deepfake detectors can be successfully attacked in both a white-box and black-box scenario [12] and that attacks prove robust to compression codecs [13], making them a very concrete threat.

However, these works limit their analyses to CNN-based detectors neglecting architectures based on transformers. Here we want to fill this gap and study the behavior of CLIP-based detectors of synthetic images in the presence of adversarial attacks and compare it with that of CNN-based detectors. It is worth to note that in computer vision there has been significant attention to analyze the adversarial robustness of

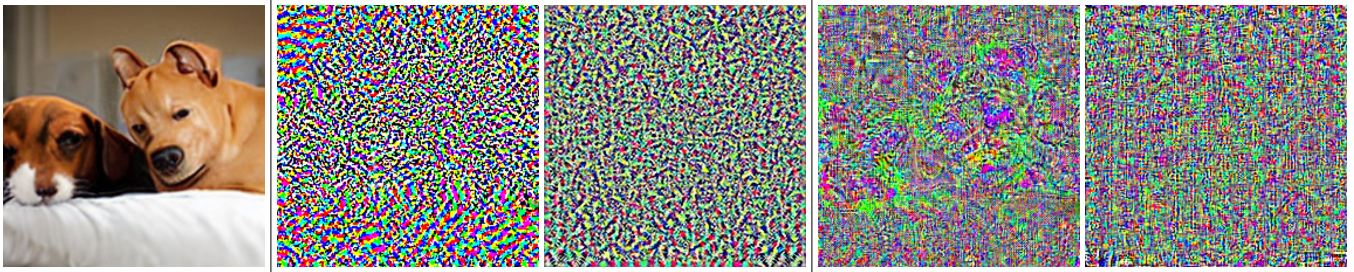


Fig. 2: From left to right: attacked image; magnified noise patterns generated by RWA and UA for a CNN-based detector; magnified noise patterns generated by RWA and UA for a CLIP-based detector. Even for these attacks we can make similar observations as done for l_2 -PGD in Fig. 1: CLIP-based attacks are more structured and show clear regular patterns.

ViT compared to CNN based solutions for the generic task of image classification [6], [14]–[16], while there is a lack of such analysis in the forensics field. In our analysis we will give a special attention to the key issue of attack transferability. Indeed, fooling a known detector is little more than a classroom exercise: the real goal is to fool all (or most) of them. In addition to presenting numerical results, we will conduct a careful analysis of the adversarial noise patterns that emerge in this process, both in the spatial and frequency domains (see Fig. 1 and Fig. 2). The results will provide interesting clues about how these detectors work and suggest possible research directions for their further improvement in terms of performance or robustness.

II. PRELIMINARIES

A. Forensic detectors

In our experiments, we consider four CNN-based detectors and four CLIP/ViT-based detectors, summarized in Table I. To conduct a fair comparison, all detectors are trained on the dataset proposed in [3] comprising real images from the COCO and LSUN datasets and synthetic images generated by Latent Diffusion Models (LDMs) [17]. In the CNN-based detector family, we first consider the popular detector proposed in [18], a simple ResNet50 pre-trained on ImageNet and optimized for synthetic image detection with compression-based augmentation and blurring. The second detector [19] differs from the first one only for one architectural change, the removal of subsampling in the first layer of the network as suggested in [19]. The same architecture is used in the third detector [3] which adopts stronger live data augmentation to improve robustness, including blurring, scaling, CutOut, compression, noise addition, and color jittering. For the last CNN-based detector, we consider a newer architecture, ConvNext Tiny [20], comparable to ResNet50 in number of parameters, also modified to avoid downsampling in the first layer.

The first two transformer-based detectors are based on the strategy proposed in [2] where features are extracted from a large pre-trained model, CLIP/ViT-L [21] in the first case, and the larger EVA-CLIP/ViT-g [22] in the second case, followed by a single-layer neural network fine-tuned for synthetic image detection. The last two detectors adopt an architecture

	Family	Network	Fine-tuning	Aug.
(1)	CNN	ResNet50	e2e	*
(2)	CNN	ResNet50 [†]	e2e	*
(3)	CNN	ResNet50 [†]	e2e	**
(4)	CNN	ConvNeXt Tiny [†]	e2e	**
(5)	CLIP	ViT-L + 1 FC layer	F	*
(6)	CLIP	ViT-g + 1 FC layer	F	*
(7)	CLIP	ViT-B + 2 FC layers	e2e	*
(8)	CLIP	ViT-B + 2 FC layers	e2e	**

(*) compression and blurring

(**) compression, blurring, scaling, cut-out, noise addition, jittering

TABLE I: List of detectors used in our experiments. All methods were trained using the dataset proposed in [3]. ViT-based networks are followed by 1 or 2 fully connected (FC) layers. Models with [†] were modified to not perform downsampling in the first layer. Fine-tuning is performed end-to-end on the entire network (e2e) or only on the last linear layer, keeping the backbone frozen (F).

consisting of CLIP/ViT-B followed by two linear layers and the whole network is fine-tuned for synthetic image detection using two different augmentation strategies proposed in [18] and [3], respectively. In Table I, we summarize the main features of such detectors.

B. Adversarial attacks

Neural networks are known to be vulnerable to adversarial attacks. By adding subtle adversarial noise to the input image, it is possible to fool the target classifier into predicting a wrong label. Forensic detectors are no exception as shown in the literature for applications as diverse as image forgery detection and source attribution [23]. Similarly, several papers [12], [24], [25] have shown that deepfake detectors can be easily attacked in a white-box scenario, i.e., when all the details of the detector are perfectly known. Even if the attacker has only partial knowledge of the detector, an attack designed for a surrogate model can be used with good success rates [12], [13]. Conversely, achieving attack transferability in a zero-knowledge scenario is non-trivial, especially if the surrogate

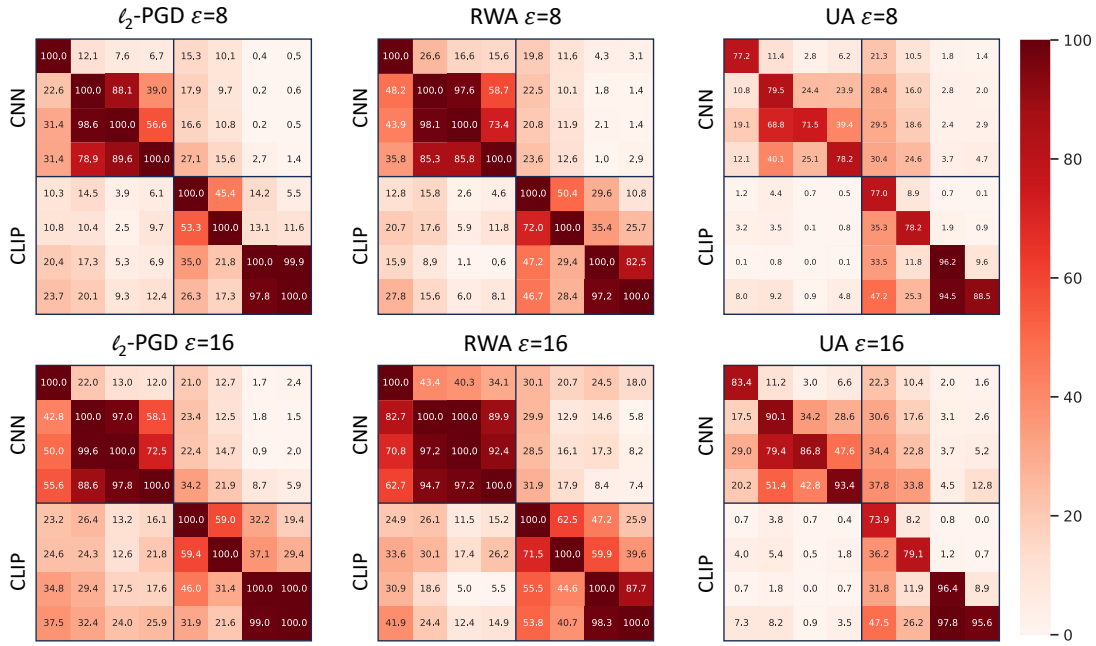


Fig. 3: Successful Attack Rate of three attacks (PGD, RWA, UA) at two strength levels ($\epsilon=8$, $\epsilon=16$) on the eight detectors of Tab. I. Cells on the diagonal correspond to white-box attacks. Off-diagonal cells correspond to transferred attacks.

and target architectures are significantly different or trained with very different protocols [26], [27].

Here, we focus on white-box attacks since our aim is to investigate the fundamental properties of different families of detectors in fully controlled conditions. However, we will also investigate in depth attack transferability, as it pertains the scenario of highest practical interest.

When the detector is perfectly known together with its parameters, θ , the attacker can compute all possible gradients of interest. This is exploited in the Fast Gradient Sign Method (FGSM) proposed originally in [28]. By computing the gradient of the model loss $L(\theta, x, y)$ with respect to the input image x one finds the direction that maximizes the disruptive effect of the attack for a given input perturbation. Following this idea, in FGSM the adversarial sample is computed as

$$x_a = x + \epsilon \text{sign}[\nabla_x L(\theta, x, y)]$$

Note that, to enforce a strict bound on the image distortion, only the sign of the gradient is taken, with the parameter ϵ controlling the strength of the perturbation.¹ In this work we will consider three different attacks, all leveraging in different ways the gradient-based approach of FGSM but with much improved performance.

The **Projected Gradient Descent (PGD)** [29] can be seen as an iterative version of FGSM. Indeed, the single iteration of FGSM is optimal only when $\epsilon \rightarrow 0$, while it may be largely sub-optimal for stronger attacks. So, PGD builds the adversarial example starting from $x^0 = x$ and updating

it iteratively as $x^{t+1} = \Pi_{B(x)} \{x^t + \alpha \mathbf{N}[\nabla_{x^t} L(\theta, x^t, y)]\}$ where $\alpha \ll 1$, $\mathbf{N}[\cdot]$ indicates normalization with reference to the adopted norm, and $\Pi_{B(x)}$ projects the result on a ball defined by the norm, centered on x , and with radius ϵ . PGD is much more effective than FGSM but still computationally affordable. Moreover, in [29] it is claimed to ensure good transferability to other architectures. In any case, it is a de-facto baseline for all studies on adversarial attacks.

The **Robust White-Box Attack (RWA)** [13] is structurally similar to PGD with some important changes meant to address forensic applications. First of all, it is explicitly developed for deepfake detection, hence a native two-class problem. More important, it takes into account the fact that, in this field, images are usually subject to various kinds of impairing transforms (compression, resizing, blurring) before being analyzed by the detector. Therefore, a generalized loss function is defined (computed in practice by sample averaging) which includes transformed versions of the original image to gain robustness to this adverse scenario.

Finally, we include the **Universal Attack (UA)** proposed in [30]. Here, the goal is not to attack a given image but rather to design a single adversarial noise sample, say Δx , such that the classifiers chooses a wrong label, $y(x + \Delta x) \neq y(x)$, for most input images. The desired perturbation is obtained by accumulating the gradients computed on a large sample of input images. Experiments show UA to be surprisingly effective, although less than targeted attacks. We include it in our analysis because it fits the classifier as a whole, hence we hope that it helps shading light on specific properties of the detectors under investigations.

¹The same parameter ϵ is used also in other attacks to control the distortion introduced in the image and is commonly taken to indicate the attack strength.

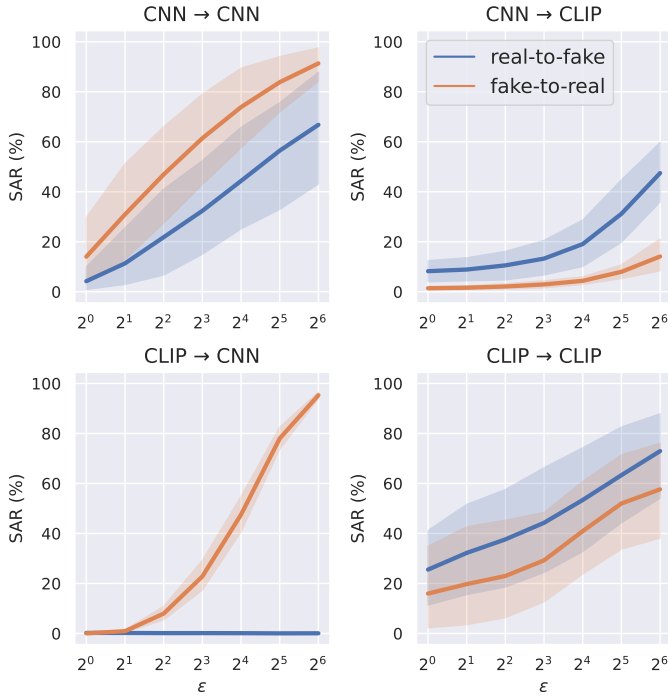


Fig. 4: Transferability of l_2 -PGD attacks as a function of the attack strength ϵ . Solid lines represent the average SAR and colored bands its standard deviation. We consider target and source detectors belonging to the same family (top-left, bottom-right) or different families (top-right, bottom-left).

III. ATTACK TRANSFERABILITY

The attacks are evaluated in terms of Success Attack Rate (SAR) measured on 1000 real images of the COCO dataset and 1000 images generated by Latent Diffusion Models (LDMs). To ensure compatibility with all detectors, images are central cropped to 224×224 pixels. In Fig. 3 we report synthetic numerical results extracted from our experimental campaign. The figure displays six matrices, corresponding from left to right to the three attacks (PGD, RWA, UA) and from top to bottom to two attack strengths ($\epsilon=8$ lighter, and $\epsilon=16$ stronger), corresponding to an average PSNR larger than 37dB and 32dB, respectively, which guarantee to avoid visual artifacts. In each matrix, the cell (i, j) reports the average success rate (SAR) observed on the j -th detector (target) using adversarial samples designed on the i -th detector (source), where $i, j \in \{1, \dots, 8\}$ span the eight detectors listed in Tab. I, 4 CNN-based and 4 CLIP-based. Cells on the diagonal show the white-box SAR for each detector. The image-targeted attacks, PGD and RWA, have a uniform success rate of 100% on all detectors and both strengths. The universal attack is somewhat less effective, as expected, but the SAR is always over 70%. These results were all largely expected.

More interesting are the results on attack transferability, given by the off-diagonal cells. We now focus on the upper-left matrix, corresponding to PGD@ $\epsilon=8$, since the other matrices show only minor differences. The color coding of the cells

(darker as the attack gets more successful) allows one to catch the big picture at a glance. Contrary to some claims in the literature, attacks are not easily transferable. To be more precise, the SAR is always very low whenever source and target models belong to different families, CNN vs. CLIP. This makes perfect sense, considering the profound architectural difference that characterize these models. However, mixed results are observed also within the same family. For example, attacks seems to transfer very well between detectors 2, 3 and 4, but not from this group to detector 1. This is very interesting, considering that detector 4 uses a ConvNeXt backbone while all the others use ResNet50, suggesting that the backbone does not impact as much on the detector behavior as other architectural choices, like removing subsampling from the first layer. In hindsight, this analysis provides valuable information on what counts for detection and what is irrelevant.

With the help of Fig. 4 we investigate attack transferability in some more detail, analyzing for l_2 -PGD separately attacks to real and synthetic images and considering a wide range of attack strengths, from $\epsilon=1$ to $\epsilon=64$. Each chart reports the average SAR obtained with source detector in one family, CNN or CLIP, and target detector in another family. Of course, even when the two families coincide, only different source and target detectors are considered. The top-left and bottom-right charts are for same-family cases, CNN \rightarrow CNN and CLIP \rightarrow CLIP. Results are as expected, with a pretty good transferability that increases smoothly with attack strength. The other two charts are more interesting. In particular, they show clearly that CNN detectors are vulnerable to fake-to-real attacks, the most dangerous for real-world cases. On the contrary, CLIP detectors are weaker with respect to real-to-fake attacks, but only at high strengths, when image distortion cannot be neglected anymore. This behavior will be more easily explained by considering the Fourier-domain analyses of next Section.

IV. FREQUENCY-DOMAIN ANALYSIS

A. Power spectra of adversarial noise patterns

It is well known that image generators cannot perfectly mimic the frequency domain behavior of natural images, sometimes introducing obvious artifacts. As a result, forensic detectors also tend to focus on specific regions of the spectrum. Therefore, by analyzing the Fourier spectrum of adversarial noise, we can gain insight into how different detectors analyze the test images [31].

In Fig. 5, for each of the three considered attacks, l_2 -PGD, RWA and UA, from top to bottom, we show the average power spectra of the adversarial noises generated to attack the eight detectors under test. In more detail, for each attack and source detector, we generate 2000 adversarial noise patterns, 1000 for real images and 1000 for fake ones, compute the squared modulus of their Fourier transform and average them to obtain the desired power spectrum.

The analysis reveals several interesting facts, first of all the differences between CNN-based and CLIP-based detectors.

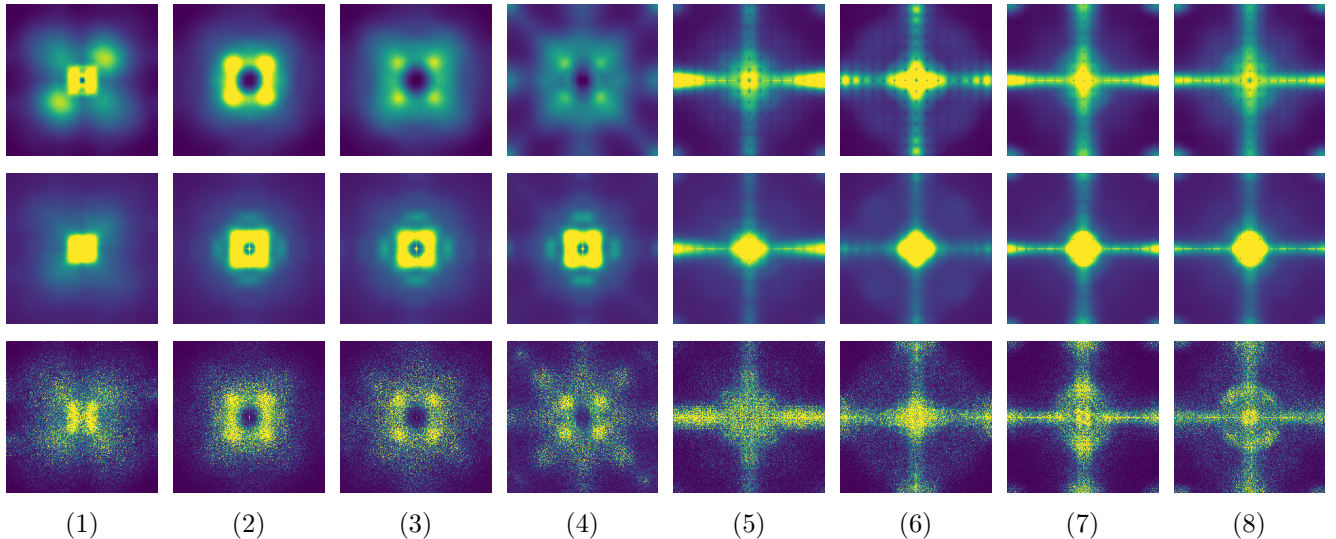


Fig. 5: Power spectra of adversarial noise patterns generated by a specific attack (rows) on a selected detector (columns). From top to bottom: l_2 -PGD, RWA, UA attacks. From left to right: detectors (1) to (8) listed in Table I.

CNN spectra² exhibit strong components at medium-high frequencies while CLIP spectra are concentrated at medium-low frequencies. This is clearly related to the backbones since, as already noted in the literature [32], CNNs are better than Transformers in capturing fine image details (high frequencies) and, conversely, the latter are better equipped to see non-local or long-range dependencies (low frequencies). In addition, CNN spectra are approximately isotropic, while CLIP spectra are markedly cross-shaped. This is due to the block processing (blocks of 14×14 or 16×16 pixels) performed by ViT in the initial stages, as also confirmed by the sinc-like behavior of the horizontal and vertical cuts of the spectra, with regular zeros, reminiscent of the transformation of a spatial rectangular window. These large differences help explain the different behavior of CNN-based and CLIP-based detectors. For example, in [33] it was observed that CLIP-based detectors are largely immune to image autoencoding-based recycling attacks, which are very effective with CNN-based detectors. Considering that such attacks mainly target the high frequencies of the signal, this fact is no longer a mystery.

On the other hand, the specific properties of the individual detectors are also important, as are the specific characteristics of the individual attacks. For example, the spectra of detector (1), in contrast to other CNN-based detectors, are more concentrated at low frequencies. We explain this by the presence of early subsampling in the detector architecture, which significantly attenuates high frequencies from the image. The augmentation protocol also plays a role. Detectors trained with strong augmentation, such as (3) and (8), consider a wider range of frequencies than their weakly augmented versions, (2) and (7). Finally, note that the spectra obtained for the RWA attack are also more concentrated at low frequencies. In our

interpretation, this is to withstand post-processing actions that tend to erase fine details.

B. On the complementarity conjecture

Before concluding this work, we want to further elaborate on the most relevant concept that emerged from the previous analysis, namely the apparent complementarity between CNN-based and CLIP-based detectors. We have already hypothesized this fact in [33]. The analysis of adversarial noise patterns has further strengthened this conjecture, which we also found confirmed in [32] where it is stated that “convolution and multi-head self-attention show opposite behaviors: the former amplifies high-frequency components while the latter reduces them”. In Fig. 6, we report the results of a further experiment along this line. The plots show the performance in terms of Accuracy (left) and Area under the ROC curve (right) of two comparable CNN-based (3) and CLIP-based (8) synthetic image detectors, as a function of the bandwidth B of a low-pass filter used to remove high-frequency components. The CLIP-based detector achieves its best performance at $B=0.2$, showing that it does not need the highest image frequencies to make a correct decision. The CNN-based detector, on the other hand, achieves a similar level only at $B=0.3$, relying largely on the 0.2-0.3 band. Therefore, they appear to focus on different portions of the frequency spectrum, thus confirming the complementarity conjecture.

V. DISCUSSION

We have studied the adversarial robustness of CLIP-based detectors for AI-generated images, also in comparison to the most popular CNN-based detectors. Analysis of the numerical results and visual inspection of the spectra of adversarial patterns allow us to draw some significant lessons, summarized below:

²Short for “power spectra of adversarial patterns generated to attack CNN-based detectors”.

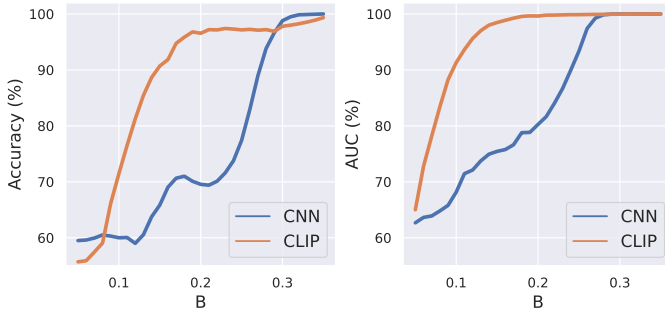


Fig. 6: Performance of a CNN-based detector and a CLIP-based detector on low-pass filtered images as a function of the filter bandwidth B . Left: Accuracy, right: AUC.

- CLIP-based detectors rely mainly on low image frequencies, in contrast to CNN-based ones that rely more on medium-high frequencies. Although our observations are specific to media forensics, it is worth noting that similar results emerged in other computer vision fields as well [15], [31];
- CLIP-based detectors are not more robust to adversarial attacks than CNN-based detectors; however, we observed very limited transferability of attacks, especially between detectors using significantly different architectures, as CNN and Transformer backbones. This finding is also in-line with recent work in image classification [16];
- The properties of adversarial pattern also depend on architectural and procedural details, but are rather stable for CLIP-based detectors;
- CLIP-based detectors are more robust than CNNs to fake-to-real attacks, which are probably the most relevant threat in a realistic scenario.

Studying robustness to adversarial perturbations is important in itself for many real-world applications. However, this study also sheds light on how detectors work and is therefore a valuable tool for developing more effective and robust detectors. Along these lines, in future work we will study the impact of adversarial attacks on generic forensic traces in AI-generated images.

ACKNOWLEDGMENT

We gratefully acknowledge the support of this research by a TUM-IAS Hans Fischer Senior Fellowship and a Google Gift. In addition, this work has received funding by the European Union under the Horizon Europe vera.ai project, Grant Agreement number 101070093, and was partially supported by SERICS (PE00000014) under the MUR National Recovery and Resilience Plan, funded by the European Union - NextGenerationEU.

REFERENCES

- [1] D. Tariang, R. Corvi, D. Cozzolino, G. Poggi, K. Nagano, and L. Verdoliva, "Synthetic Image Verification in the Era of Generative AI: What Works and What Isn't There Yet," *IEEE Security & Privacy*, vol. 22, pp. 37–49, 2024.
- [2] U. Ojha, Y. Li, and Y. J. Lee, "Towards universal fake image detectors that generalize across generative models," in *CVPR*, 2023, pp. 24 480–24 489.
- [3] 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, pp. 1–5.
- [4] Z. Sha, Z. Li, N. Yu, and Y. Zhang, "DE-FAKE: Detection and Attribution of Fake Images Generated by Text-to-Image Generation Models," in *ACM SIGSAC Conference on Computer and Communications Security*, 2023, pp. 3418–3432.
- [5] R. Amoroso, D. Morelli, M. Cornia, L. Baraldi, A. Del Bimbo, and R. Cucchiara, "Parents and Children: Distinguishing Multimodal DeepFakes from Natural Images," *ACM Transactions on Multimedia Computing, Communications and Applications*, may 2024.
- [6] Y. Bai, J. Mei, A. Yuille, and C. Xie, "Are Transformers More Robust Than CNNs?" in *NeurIPS*, 2021, pp. 26 831–26 843.
- [7] M. Barni, M. C. Stamm, and B. Tondi, "Adversarial Multimedia Forensics: Overview and Challenges Ahead," in *EUSIPCO*, 2018, pp. 962–966.
- [8] L. Lin, N. Gupta, Y. Zhang, H. Ren, C.-H. Liu, F. Ding, X. Wang, X. Li, L. Verdoliva, and S. Hu, "Detecting multimedia generated by large AI models: A survey," *arXiv preprint arXiv:2402.00045*, 2024.
- [9] F. Marra, D. Gragnaniello, L. Verdoliva, and G. Poggi, "Do GANs Leave Artificial Fingerprints?" in *MIPR*, 2019, pp. 506–511.
- [10] N. Yu, L. Davis, and M. Fritz, "Attributing Fake Images to GANs: Learning and Analyzing GAN Fingerprints," in *ICCV*, 2019, pp. 7556–7566.
- [11] R. Corvi, D. Cozzolino, G. Poggi, K. Nagano, and L. Verdoliva, "Intriguing properties of synthetic images: from generative adversarial networks to diffusion models," in *CVPR Workshop*, 2023, pp. 973–982.
- [12] N. Carlini and H. Farid, "Evading deepfake-image detectors with white- and black-box attacks," in *CVPR Workshop*, 2020, pp. 658–659.
- [13] S. Hussain, P. Neekhara, M. Jere, F. Koushanfar, and J. McAuley, "Adversarial Deepfakes: Evaluating Vulnerability of Deepfake Detectors to Adversarial Examples," in *WACV*, 2021, pp. 3348–3357.
- [14] S. Bhojanapalli, A. Chakrabarti, D. Glasner, D. Li, T. Unterthiner, and A. Veit, "Understanding Robustness of Transformers for Image Classification," in *ICCV*, 2021, pp. 10 231–10 241.
- [15] P. Benz, S. Ham, C. Zhang, A. Karjauv, and I. S. Kweon, "Adversarial Robustness Comparison of Vision Transformer and MLP-Mixer to CNNs," in *BMVC*, 2021.
- [16] K. Mahmood, R. Mahmood, and M. Van Dijk, "On the Robustness of Vision Transformers to Adversarial Examples," in *ICCV*, 2021, pp. 7838–7847.
- [17] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *CVPR*, 2022, pp. 10 684–10 695.
- [18] S.-Y. Wang, O. Wang, R. Zhang, A. Owens, and A. A. Efros, "CNN-generated images are surprisingly easy to spot... for now," in *CVPR*, 2020, pp. 8692–8701.
- [19] D. Gragnaniello, D. Cozzolino, F. Marra, G. Poggi, and L. Verdoliva, "Are GAN generated images easy to detect? A critical analysis of the state-of-the-art," in *ICME*, 2021, pp. 1–6.
- [20] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A ConvNet for the 2020s," in *CVPR*, 2022, pp. 11 976–11 986.
- [21] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*, 2021, pp. 8748–8763.
- [22] Q. Sun, Y. Fang, L. Wu, X. Wang, and Y. Cao, "EVA-CLIP: Improved Training Techniques for CLIP at Scale," *arXiv preprint arXiv:2303.15389*, 2023.
- [23] M. Barni, W. Li, B. Tondi, and B. Zhang, "Adversarial Examples in Image Forensics," in *Multimedia Forensics*. Springer, 2022.
- [24] P. Neekhara, B. Dolhansky, J. Bitton, and C. Ferrer, "Adversarial Threats to DeepFake Detection: A Practical Perspective," in *CVPR*, 2021, pp. 923–932.
- [25] S. Jia, C. Ma, T. Yao, B. Yin, S. Ding, and X. Yang, "Exploring Frequency Adversarial Attacks for Face Forgery Detection," in *CVPR*, 2022, pp. 4103–4112.
- [26] X. Zhao and M. C. Stamm, "Making Generated Images Hard To Spot: A Transferable Attack On Synthetic Image Detectors," in *ICPR*, 2022, pp. 70–84.

- [27] Y. Wang, X. Ding, L. Ding, R. Ward, and Z. J. Wang, "Perception Matters: Exploring Imperceptible and Transferable Anti-forensics for GAN-generated Fake Face Imagery Detection," *Pattern Recognition Letters*, vol. 146, pp. 15–22, 2021.
- [28] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and Harnessing Adversarial Examples," in *ICLR*, 2015.
- [29] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards Deep Learning Models Resistant to Adversarial Attacks," in *ICLR*, 2018.
- [30] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal Adversarial perturbations," in *CVPR*, 2017, pp. 1765–1773.
- [31] J. Bai, L. Yuan, S.-T. Xia, S. Yan, Z. Li, and W. Liu, "Improving Vision Transformers by Revisiting High-frequency Components," in *ECCV*, 2022, pp. 1–18.
- [32] N. Park and S. Kim, "How Do Vision Transformers Work?" in *ICLR*, 2022.
- [33] D. Cozzolino, G. Poggi, R. Corvi, M. Nießner, and L. Verdoliva, "Raising the Bar of AI-generated Image Detection with CLIP," in *CVPR Workshop*, 2024, pp. 4356–4366.