

WILD: a new in-the-Wild Image Linkage Dataset for synthetic image attribution

Pietro Bongini^{*}, Sara Mandelli[†], Andrea Montibeller[‡], Mirko Casu[§], Orazio Pontorno[§], Claudio Vittorio Ragaglia[§]
 Luca Zanchetta[¶], Mattia Aquilina[¶], Taiba Majid Wani[¶], Luca Guarnera[§], Benedetta Tondi^{*}
 Giulia Boato[‡], Paolo Bestagini[†], Irene Amerini[¶], Francesco De Natale[‡], Sebastiano Battiato[§], Mauro Barni^{*}

^{*} University of Siena, Department of Information Engineering and Mathematics, Italy

[†] Politecnico di Milano, Department of Electronics, Informatics, and Bioengineering, Italy

[‡] University of Trento, Department of Information Engineering and Computer Science, Italy

[§] University of Catania, Department of Mathematics and Computer Science, Italy

[¶] Sapienza University of Rome - Department of Computer, Control and Management Engineering, Italy

Abstract—Synthetic image source attribution is an open challenge, with an increasing number of image generators being released yearly. The complexity and the sheer number of available generative techniques, as well as the scarcity of high-quality open source datasets of diverse nature for this task, make training and benchmarking synthetic image source attribution models very challenging. WILD¹ is a new in-the-Wild Image Linkage Dataset designed to provide a powerful training and benchmarking tool for synthetic image attribution models. The dataset is built out of a closed set of 10 popular commercial generators, which constitutes the training base of attribution models, and an open set of 10 additional generators, simulating a real-world in-the-wild scenario. Each generator is represented by 1,000 images, for a total of 10,000 images in the closed set and 10,000 images in the open set. Half of the images are post-processed with a wide range of operators. WILD allows benchmarking attribution models in a wide range of tasks, including closed and open set identification and verification, and robust attribution with respect to post-processing and adversarial attacks. Models trained on WILD are expected to benefit from the challenging scenario represented by the dataset itself. Moreover, an assessment of seven baseline methodologies on closed and open set attribution is presented, including robustness tests with respect to post-processing.

Index Terms—Synthetic image attribution in-the-wild, multimedia forensics, deep learning, benchmarking dataset.

I. INTRODUCTION

Advancements in Generative AI (GenAI), particularly Diffusion Models (DMs) and Generative Adversarial Networks (GANs), have revolutionized multimedia content creation, producing highly realistic synthetic images that blur the line with reality [1]–[3]. These innovations, while impressive, challenge multimedia forensics by complicating the attribution of synthetic images to their source generators, a task that is both intricate and essential. In fact, synthetic image attribution (sometimes referred to as image linkage), that is, the identification of the generative model used to produce a given image, is vital for forensic analysis, intellectual property protection, and counter AI-driven misinformation and impersonation. However, current methods struggle to generalize,

especially in open-set scenarios involving unknown or new generators [4], [5].

Some datasets are available that can be used to address this task. For example, the Synthetic Faces High Quality - Text2Image (SFHQ-T2I) dataset [6] provides 120,000 high-quality synthetic faces (1024 × 1024 resolution) from models such as FLUX.1 (Pro, dev, schnell), Stable Diffusion XL, and DALL·E 3. However, this dataset presents a strongly uneven image distribution across generators and lacks prompt-image linkage to mirror real-world attribution difficulties. The ArtiFact dataset, proposed in [7], comprises 25 diverse generators, object categories, and real-world challenges. Nonetheless, ArtiFact’s image resolution is very limited (i.e. 200 × 200) and, once again, images are not uniformly distributed among generators. The GenImage dataset [8] comprises eight different generative models, each used to produce synthetic images for the 1,000 distinct labels in ImageNet [9], ensuring a nearly equal distribution across classes. In contrast, this dataset does not include any form of image processing, resembling only ideal laboratory conditions. Moreover, none of these datasets include an open set for effective simulations of what attribution models trained on them can achieve in the wild. Additionally, commercial generators, even the most popular and utilized ones, are scarce if present at all in such datasets: among the datasets cited above, the only examples are Midjourney in GenImage and DALL·E 3 in SFHQ-T2I.

The lack of datasets containing high-resolution images, uniformly distributed across synthetic generators (including commercial ones as well), with editing operations applied, and of an open set for simulating in-the-wild conditions, prevents the possibility of evaluating attribution models on realistic scenarios. To tackle these issues, we present WILD, a novel “in-the-Wild” Image Linkage Dataset tailored for synthetic image attribution. The dataset focuses on the source attribution, without taking into account real vs fake detection. The dataset is composed exclusively of head-and-shoulder images of synthetic people from fixed prompts. WILD features a closed set containing images produced by 10 popular text-

¹The dataset is available at: <https://www.kaggle.com/datasets/pietrob92/wild-in-the-wild-image-linkage-dataset>



Fig. 1: Two groups of ten images generated by the closed set generators with the same text prompt. The generators are identified by the codes described in Section II-B. The prompts used to generate the images are: (top, a-j): prompt 439, “Image of an Italian young man who is wearing a pullover. He has a pointy face and full lips. The subject is looking into the camera in happiness. The image is taken with a city in the background”; (bottom, k-t): prompt 266, “A proud British woman in her 20s with hooded brown eyes and a small nose is looking into the camera. She is wearing a sky blue fedora and a necklace. The subject is portrayed with a golf course in the background”.

to-image generators (1,000 images each), among which six commercial platforms, and an open set section, obtained by using 10 other generators contributing 1,000 images each.

In the closed set, image linkage allows to retrace the prompt each image was produced with. This is achieved by generating 1,000 prompts with an ad-hoc prompting script and using each prompt once on every generator. The open set is instead conceived as a simulation of what can happen in the wild. It is composed of generators with very different characteristics, in order to better cover the possible image features that can be encountered by the attribution models once in use. To reflect real-world conditions, half of the images have undergone three post-processing steps (e.g., rotation, resizing, compression, and many others): for these images, we release the plain version, and the post-processed versions with one, two, and three operators. The contributions of this paper are:

- We release WILD, a challenging source attribution dataset of head-and-shoulder synthetic faces. It includes a closed set produced by some of the most popular synthetic generators. The images are generated through randomized and evenly distributed prompts across generators. This linkage between images and prompts helps to minimize biases and improve the understanding of both generation and attribution mechanisms. WILD also includes a large and varied open set of synthetic human faces, simulating realistic in-the-wild scenarios.
- To better reflect images from the real-world, e.g., those shared through social media platforms, we apply several common post-processing operations. In this way, we allow testing the robustness of forensic detectors.
- We present the results of several experiments obtained by applying state-of-the-art source attribution baselines to the WILD dataset. This allows to assess the dataset characteristics and to benchmark new attribution models based on their results. The experiments include closed-set attribution, open-set attribution and robustness evaluation.

II. WILD DATASET

WILD is designed as an Image Linkage dataset for benchmarking source attribution methods, both in closed set and

open set scenarios. To simulate real-world conditions, the closed set includes some of the most popular commercial and open-source text-to-image generators. The open set, instead, contains a large and varied collection of different image generators (including both DMs and GANs), to summarize the widest possible range of generative techniques. The images of the dataset were generated with a total of 20 generators, 10 for the closed set and 10 for the open set. Each generator is represented by 1,000 images, for a total of 20,000 images. Half of the images in each subset were post-processed by one, two, or three image-processing operators, for a grand total of 50,000 images. In the following, we provide additional details on the composition of the dataset and its construction procedure.

A. Closed Set

The closed set was built using 10 of the most popular commercial and open source text-to-image generators. Each generator was used to generate 1,000 uncompressed images. To avoid significant biases, we developed an ad-hoc prompting methodology. Each prompt was used to generate 10 images, one for each generator. In this way, we ensured that every class contains the same data distributions. In addition, in this way we created a direct link between each image and its prompt. In Fig. 1, we report some examples of images generated with the same text prompt from each generator.

B. Closed set generators

The generators used to build the closed set are listed below, together with a description of the hyperparameters used. For the sake of readability, in the tables and figures of the following sections, these models are indicated by the three-character codes introduced here.

1) Adobe Firefly (ADF): this is the proprietary solution released by Adobe to generate extremely realistic images and edit existing photos using text descriptions [10]. Specifically, we selected the second version of the Firefly text-to-image generator, which allows to create 1024×1024 images.

2) *DALL·E 3 (DE3)*: this is a commercial text-to-image generator by Open-AI, ensuring consistency between the textual description and the visual output [11]. We generated 1024×1024 images with the standard quality settings.

3) *FLUX.1 (FX1)*: this is a diffusion-based text-to-image framework freely provided by Black Forest Labs, optimized for high-resolution, photorealistic outputs with minimal computational cost [12]. For dataset creation, we generated 1024×1024 images using a guidance scale of 3.5, 30 inference steps, and a 512-token max sequence length.

4) *FLUX 1.1 Pro (FPX)*: this is an advanced diffusion-based text-to-image framework by Black Forest Labs, designed for fast, high-resolution, photorealistic generation [13]. We generated 1024×1024 images, using a fixed seed of 42, a safety tolerance of 2, and without prompt upsampling.

5) *Freepik (FPK)*: this is a commercial text-to-image generative model [14]. We generated 1024×1024 images with parameters: 8 inference steps, a guidance scale of 1.0, a photo style with pastel color, warm lighting, and portrait mode.

6) *Leonardo AI (LEO)*: this is a commercial text-to-image generator [15]. The images were produced using the Kino XL model with the following parameters: a size of 1024×1024 , portrait style, a guidance scale of 7, 15 inference steps, and the “alchemy”, “photoReal” and “photoRealVersion” parameters set to “True”, “True” and “v2”, respectively.

7) *Midjourney (MDJ)*: this is a commercial text-to-image generative framework [16]. We used model version 6.1 and a size of 1024×1024 . In this case, a negative prompt was specified to avoid excessively unrealistic images².

8) *Stable Diffusion 3.5 Large (S35)*: this is an advanced text-to-image model freely provided by StabilityAI [17], [18]. We used 50 inference steps, a guidance scale of 3.0, and a 512-token max sequence length to produce 1024×1024 images.

9) *Stable Diffusion XL Turbo (SXL)*: this is another text-to-image model freely released by StabilityAI [19]. To ensure a good visual quality for the synthesized images, we used a fixed resolution of 512×512 pixels and 2 inference steps.

10) *Starry AI (STR)*: this is a commercial text-to-image generator [20]. We generated 1024×1024 images without high-resolution enhancement and using 20 inference steps.

C. Prompt Generation

The same 1,000 head-and-shoulder prompt descriptions were used for all the generators in the closed set. In this scope, it is expected that an attribution model can rely on semantic image features that might be relative to the person, clothing, background, or maybe on low-level features that persist throughout the whole image. As a consequence, it is fundamental to introduce a certain variability in the images across all these characteristics. This variability was injected through the text prompts. We generated the prompts automatically. A simple token system allowed various characteristics to be introduced in each prompt, randomizing or rotating them.

²Negative prompt used for Midjourney: *unrealistic, cartoon, anime, painting, illustration, greyscale, sepia, drawing, sketch, fantasy, elves, orcs, strange nose*.

Five different sentence structures were used to play with the different nuances of interpretation each model can have for different text structures. Prompts were generated in blocks of 100, rotating every sentence structure twice (for a total of 200 prompts per structure). The sentence structures can be summarized as follows³:

- A: <subject> with <facial features> wearing <clothing> <be> looking <looking direction> (in <expression(emotion)>) (<background>).
- B: Picture of a/an (<expression(adjective)>) <subject> with <facial features>. <pronoun> <wear> <clothing> and <look> <looking direction> (<background>).
- C: Image of <subject> who is wearing <clothing> (, <background>). The subject has <facial features> and (shows <expression(emotion)> while)/(is) looking <looking direction>.
- D: Head and shoulder picture of (<expression(adjective)>) <subject> with <facial features>. <pronoun> is looking <looking direction> (<background>) and is wearing <clothing>.
- E: Portrait of (<expression(adjective)>) <subject> with <facial features>. <pronoun> <wear> <clothing>. The subject <look> <looking direction> (<background>).

To avoid biases, we managed data distributions in two different ways: subjects were described according to their gender, together with their age category (8 categories equally distributed over 80% of the prompts) or profession (20% of the prompts). Each gender covers 50% of the prompts. These characteristics were rotated to ensure a bias-free dataset along the respective axes. For profession-based descriptions, the profession was randomly drawn from a uniform distribution of 110 professions. The pronoun token followed the subject’s gender. The verb tokens were simply conjugated according to the rest of the sentence. An ethnicity/nationality descriptor was randomly selected for each subject.

The <expression> token varies depending on the sentence structure and can describe the expression (e.g., smiling, happy) or emotion shown by the subject (e.g., joy, sadness). The emotion/expression was randomly chosen from a distribution.

The <facial features> token corresponds to exactly two features for every subject randomly drawn from a uniform distribution. Depending on the type of feature, subsequent random choices were made, including hair color, hairstyle, nose shape, eye color, or other similar details. Two instances of the same type were not allowed. To ensure correctness, some facial features were moved through the sentence from the position of the <facial features> token to the front of the <subject> token, for example: “Picture of a *blonde* <subject> with *freckles...*”.

³The parts in round brackets, describing expression and background, were not always included to increase variability. The parts in angle brackets represent characteristics that vary across prompts.



Fig. 2: Some images from the open set section of WILD (one image per generator). The generators are identified by the codes described in Section II-D.

The <clothing> token describes one, two, or three clothing items. If multiple items appear, they must be from different categories: main clothing, headwear, accessories. Each category has a given probability of appearing, and when no category appears, the random selection is repeated until at least one does. The specific item to include in the description is selected from a uniform distribution for each category. Depending on the item, a following randomization can be introduced to select its pattern, color, or material (e.g., respectively: a plaid scarf, a yellow t-shirt, a silver necklace).

The <looking direction> token introduced another random choice on where the subject is looking. This corresponded to “into the camera” with a probability of 0.8, or to “up”, “down”, “left”, “right” with a probability of 0.05 each.

The <background> token, when present, was randomized between three different description modes: color, pattern, scenery. If color was selected, a random color was specified (e.g. “... on a blue background.”). If pattern was selected, the prompt describes a pattern, item or group of items that should appear behind the subject (e.g. “... with leaves in the background” or “... with a fountain in the background”). If scenery was selected, the subject is described as located in a particular place, on which the background generation will depend (e.g. “... in a village”, “... on the grass”, “... in the mountains”).

The generated prompts were validated through a block of filters to avoid bad prompts and absurd co-occurrences of characteristics (e.g., a child with a beard). The prompts were then used to generate images. Each prompt that failed the validation or did not result in a successful generation with all the ten closed set generators was substituted with a new one, using the same sentence structure, gender, and age category (or profession), but randomizing again all the other characteristics.

Images with the same prompt are expected to have similar semantic characteristics, yet some differences are introduced by the different generative methods. These differences are due to model architecture, filters, training base (and biases), and parameters, and they are fundamental for source attribution, together with the low-level model signature. An example of the ten images generated with a prompt from our dataset is shown in Fig. 1 (top row, a-j), where it is evident how models can give different interpretations to prompt features: the pullover, the background, and the facial features are all interpreted with different nuances, also from a semantic point of view. A similar example is shown in the bottom row (k-t) of Fig. 1. In this case, it can be observed how prompt specifications are slightly (or sometimes completely) misinterpreted by generators: the fedora specified in the prompt is not always a fedora, but

maybe another type of hat, the necklace is ignored by some models. Moreover, some details are added or extended to cover what is not specified in the prompt: the rest of the clothing might take the same color specified for the fedora or can have a different pattern, like those produced by Leonardo AI and Stable Diffusion 3.5. When a background is not specified, instead, a plausible one might be inferred from other details.

D. Open Set

To simulate a real-world case scenario where detection models need to attribute generators in the wild, WILD features an open set section with 10,000 synthetic head and shoulder images obtained from another group of 10 generators. In this case, we set no constraints on the prompts and data distributions. Moreover, to create the open set, we did not use exclusively text-to-image models. In the following, we introduce the list of generative methods used for the open set. Also in this case, we use a three-character code to identify them in tables and figures: DALL-E 1 (DE1) [21], Deep AI (DPA) [22], Hotpot AI (HPA) [23], NVIDIA Sana PAG (NVS) [24], Stable Cascade (SCD) [25], Stable Diffusion Attend&Excite (SAE) [26], StyleGAN (SG1) [27], StyleGAN2 (SG2) [28], StyleGAN3 (SG3) [29], Tencent Hunyuan (THY) [30]. Some examples of images included in the open set are displayed in Fig. 2, showing one image for each generator of the open set.

E. Post-Processing

For post-processing, we applied various image enhancement and processing techniques to half of the images of the dataset (both closed set and open set). This simulates real-world conditions where manipulated images may undergo transformations that obscure the detection cues. Each image underwent one to three randomly selected post-processing operations chosen from a set of ten transformations: JPEG/WebP compression with a quality factor in the range [50, 100]; random central cropping with up to 10% border loss; resizing with a scaling factor in the range [0.4, 2]; rotation by $\pm 0.5^\circ$ with aspect-preserving cropping; contrast/brightness adjustment with factors in the range [0.8, 1.2]; Gaussian blur with a radius in the range [0.5, 1.5]; grayscale conversion; super-resolution using pre-trained models from NinaSR [31], CARN [32], EDSR [33], RDN [34], or RCAN [35]; JPEG AI compression [36] with variable Bit Per Pixel (BPP) compression rate in the set [0.5, 0.75, 1, 1.25, 1.50]. Post-processed images were organized by their transformation depth (1 to 3 steps), enabling controlled testing of detector robustness against increasingly complex modification chains. The post-processed images were not included in the training set and were used exclusively during testing.

F. Dataset Split

The closed set is released with a predefined split to facilitate cross-comparisons. The split was operated at the prompt level, to avoid images with the same prompt ending up in different sets, and to ensure a perfect class balance between closed set generators. The training, validation, and test sets contain respectively 5,000 images (500 prompts), 2,000 images (200 prompts), and 3,000 images (300 prompts). The splitting procedure took into account the prompt blocks described in Section II-C and respected the global data distributions, keeping biases at a minimum level. Each set has the same percentage of prompts for every gender (50%) and for every age category and the same percentage of profession-based prompts. Moreover, the five prompt structures are distributed equally among the sets. In contrast, the distributions of ethnicity/nationality, facial features, clothing, background, and expression/emotion of the characters were randomized during the prompt building procedure. As a consequence, we did not enforce their uniform distribution among sets. All the test set images have post-processed versions. The open set instead was not split, as it is not intended to be used for training.

III. BASELINES AND RESULTS

In the following, we assess the performance of some common attribution baselines on our proposed WILD dataset. We start by describing the baselines selected for the benchmark, then we show their results on closed set attribution and open set attribution.

A. Baseline methods

1) *Clip Feature Classifiers*: CLIP (Contrastive Language–Image Pre-training) [37] is used as a feature extractor in a wide range of tasks. Being pre-trained on large amounts of data, it allows to train classifiers directly on its features, even when few data are available. We used CLIP to extract image features without specifying the prompt. The features were then passed to a small classification head, which was either a Multi-Layer Perceptron (MLP) or a Support Vector Machine (SVM). We used CLIP Large with input size 336. The hyperparameters of the two models were tuned on the validation set. For CLIP+MLP, we used an MLP with two hidden layers of 512 and 256 units, respectively, both with ReLU activation and softmax output layer. The MLP was trained with Adam [38], an initial learning rate of 10^{-3} and $l_{2,\alpha}$ of 10^{-4} , for a maximum of 1,000 epochs and the validation loss as early stopping criterion, with patience 10 (restoring the best weights). The SVM instead used a Radial Basis Function (RBF) kernel, with C equal to 1.0, a tolerance of 10^{-3} , degree 3, and no maximum number of iterations.

2) *DE-FAKE*: In [39], Sha et al. introduced a novel framework for synthetic image source attribution: a hybrid classifier for robust detection and attribution across different generative techniques. The model exploits multimodal features extracted from both the input image and its prompt description through CLIP’s image and text encoder [37]. When prompt descriptions are not available, they are estimated using Blip2

[40]. During our experiments, DE-FAKE [39] was trained and tested following the procedures and using the hyperparameters proposed in the original publication.

3) *Standard Convolutional Neural Networks (CNNs)*: As additional baselines, we considered classic CNNs which have been successfully used in many forensic tasks. Specifically, we selected the EfficientNetB4 [41], XceptionNet [42], and ResNet50 [43] architectures. We trained them following a similar approach to that proposed in [44]: each CNN works at patch level, considering square RGB patches of 96×96 pixels as input and providing a single score per patch. At deployment stage, given an image, we randomly extracted 50 patches and computed the final score for the image as the arithmetic mean of the patch scores. We adopted this approach as it demonstrated strong robustness to compression and resizing operations, thanks to an extensive set of augmentations incorporated into the training process [44]. The ability to extract small patches from the query image reduces the dependency on its semantic content, allowing for a sharper focus on the synthetic generation artifacts.

4) *Vision Transformer Classifier*: Vision Transformers (ViT) [45] are a valid alternative to CNNs for synthetic image detection. ViTs leverage self-attention to capture long-range dependencies. Along with the other ViT-based methods, which exploit the ViT for feature extraction, we also employed a ViT directly as a classifier. For the sake of clarity, we call it Vision Transformer Classifier (VTC) in the following. We employed the ViT-Base model, pre-trained on ImageNet [9] and ImageNet-21k [46] and fine-tuned it for synthetic image attribution on our closed set. Images were split into 16×16 patches, embedded, and processed by the transformer encoder. We aggregated patch-level predictions to obtain the final classification score, enhancing resilience against localized distortions.

B. Closed Set Attribution

We evaluated the models’ performance in terms of balanced accuracy for closed-set attribution. This is the simplest task, as the evaluation is limited to the ten generators encountered during training. The results in Table I show that CNN-based methods, known to learn low-level features, significantly outperform other baselines on plain images. In contrast, their performance drops sharply when evaluated on post-processed images, reporting the worst results among all the investigated methodologies. The VTC method, which relies on higher-level features, proves to be more robust to post-processing. These findings align with recent observations in [47], which indicate that ViT-based solutions are more suitable for use in the wild.

C. Open Set Results

For the open set attribution task, we adopted two different approaches. We started by modeling open set attribution as a binary classification problem, i.e., separating closed set samples from open set ones. Then, we tested the baselines on multi-class attribution, considering, together with the open set samples, all the closed set classes.

TABLE I: Balanced accuracy on the closed set attribution task for all the baselines, on the plain and post-processed (1,2,3 steps) image versions.

Baseline	Plain	1 Step	2 Steps	3 Steps
CLIP+MLP	96.67%	87.13%	80.13%	72.40%
CLIP+SVM	95.30%	86.50%	78.93%	72.63%
DE-FAKE	94.00%	87.00%	83.00%	78.00%
EfficientNetB4	100.0%	84.43%	73.03%	59.93%
XceptionNet	100.0%	84.00%	72.67%	58.70%
ResNet50	99.90%	84.87%	73.53%	60.63%
VTC	95.80%	93.40%	91.50%	88.67%

TABLE II: Open-set vs closed-set classification results (without post-processing).

Baseline	ROC-AUC	EER	TPR@FPR=0.05
CLIP+MLP	0.876	0.202	0.373
CLIP+SVM	0.872	0.210	0.385
DE-FAKE	0.836	0.244	0.485
EfficientNetB4	0.997	0.028	0.984
XceptionNet	0.992	0.046	0.961
ResNet50	0.985	0.061	0.923
VTC	0.890	0.181	0.374

1) *Binary classification:* We used the scores provided by the attribution models over the samples coming from the test set of the closed set and from the open set. We selected the maximum softmax scores achieved over all the ten closed set classes. Then, we exploited the distribution of these scores to classify the test samples as coming from the “closed set” class (positive) or from the “open set” one (negative). The “open set” class is reasonably associated with low softmax probabilities. To evaluate the models’ performance, we measured the binary ROC-AUC, the Equal Error Rate (EER), and the True Positive Rate evaluated at False Positive Rate equal to 5% (TPR@FPR=0.05). Table II reports the results on non-post-processed test set samples. Similarly to what has been observed for the closed-set, these results highlight the strong capability of the three CNN-based methods to discriminate closed-set samples from open-set ones, with EfficientNetB4 achieving the best results on all the metrics.

In Table III, we compare the three best state-of-the-art methodologies (considering only a single approach among the CNN-based ones) over post-processed images. As in the closed-set case, VTC tends to be the most robust approach, even if the TPR@FPR=0.05 is still quite low and sometimes outperformed by CNN-based solutions.

Finally, Fig. 3 presents the ROC curves obtained for both plain and post-processed images, confirming our previous findings. CNN-based methods significantly outperform other baselines on plain images, with EfficientNetB4 demonstrating the best ROC curve, closely followed by ResNet50 and XceptionNet. However, their advantage diminishes on post-processed images. With three post-processing steps, VTC emerges as the best-performing model, while the other methods exhibit comparable performance. This again highlights that high-level features like those extracted by ViT-based methods contribute to greater model robustness.

2) *Multi-class classification results:* In this setup, we evaluated the performance of the detectors in attributing both the

TABLE III: Binary closed vs open set classification results for CLIP+MLP, EfficientNetB4, and VTC on post-processed images.

Metric(Steps)	CLIP+MLP	EfficientNetB4	VTC
ROC-AUC (1)	0.781	0.872	0.857
ROC-AUC (2)	0.720	0.763	0.829
ROC-AUC (3)	0.662	0.659	0.801
EER (1)	0.284	0.202	0.215
EER (2)	0.341	0.308	0.240
EER (3)	0.384	0.389	0.271
TPR@FPR.05 (1)	0.180	0.560	0.256
TPR@FPR.05 (2)	0.122	0.305	0.197
TPR@FPR.05 (3)	0.089	0.163	0.169

closed set and open set samples. To do this, we followed a two-step approach. First, we fixed a rejection threshold on the softmax probabilities. Then, if no class probability was larger than the threshold, the example was rejected, and classified as unknown or “open set”, otherwise we attributed it to the class with maximum probability. We fixed the rejection threshold on the validation set samples. We recall that the validation set contains only non-processed, closed set samples. We then extracted the maximum softmax scores and set the threshold to obtain a rate of incorrect detections (i.e., samples detected as being “open set”) equal to 5%. At test time, we applied the rejection threshold to samples of the closed set (test set) and the open set. We counted how many samples were correctly rejected but also how much the rejection impacted the attribution capabilities on known classes. Fig. 4 shows the obtained confusion matrices. Notice that the “open set” has been added as an additional class. For the sake of readability, we do not report the results of all the baselines, but only the most interesting ones. The Correct Classification Rate (CCR) of every closed set class can be read in the diagonal cells, as well as the Correct Rejection Rate (CRR) of the open set (last diagonal cell). In this task, XceptionNet neatly outperforms all the other baselines: having a perfect CCR of 1.0 on all closed set classes (which is true also for EfficientNetB4 and ResNet50), it also shows reasonably good rejection performance with CRR equal to 0.7.

Concerning robustness against post-processing operations, we report in Table IV the average CCR and CRR of the three best methodologies. Concerning the closed set performance, it is interesting to notice that XceptionNet degrades rapidly its CCR on every class, confirming our previous considerations on the poor robustness (regarding the closed set attribution capabilities) of CNN-based baselines to editing operations. A similar behavior, even if less steep, is shown by DE-FAKE, whose CCR decreases as the number of editing steps increases. The best closed set attribution method appears again to be VTC, with a relatively large CCR even in the presence of post-processing. With regard to the *open set*, XceptionNet proves to be the best methodology overall, even if all the investigated techniques show a drop in performance when the images go through multiple processing steps.

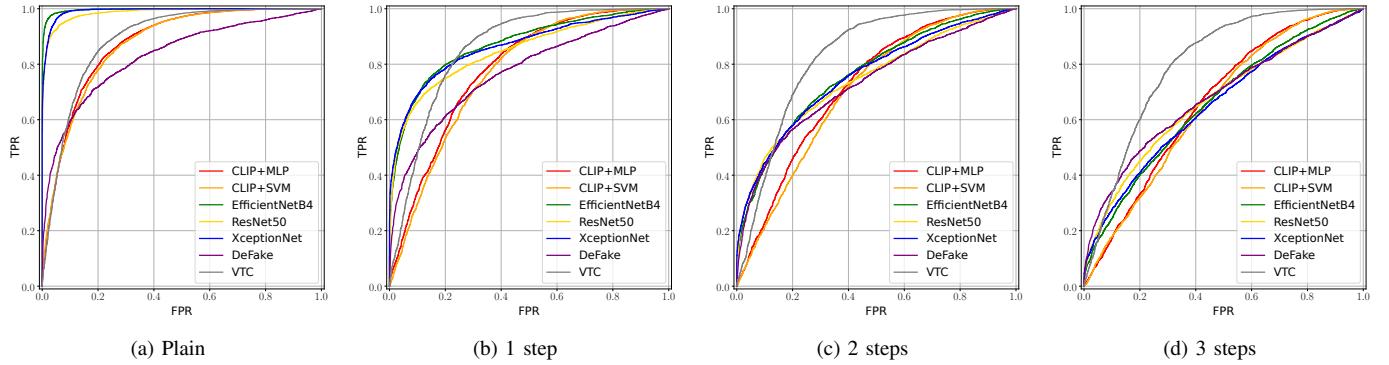


Fig. 3: ROC Curves of baselines in the binary classification task between open set (positive class) and closed set (negative task), with different levels of image post-processing: (a) plain images, (b) 1 post-processing step, (c) 2 post-processing steps, (d) 3 post-processing steps.

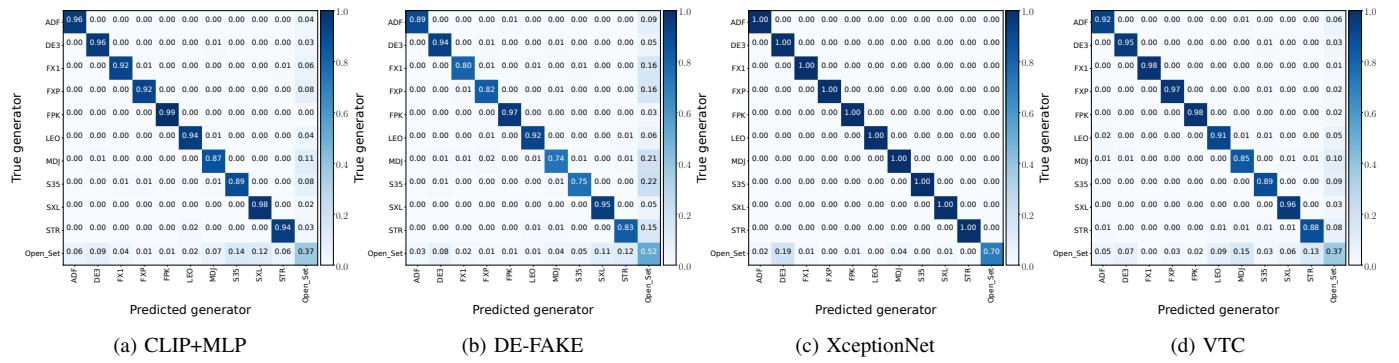


Fig. 4: Confusion matrices obtained by CLIP+MLP (a), DE-FAKE (b), XceptionNet (c), and VTC (d) on the open set attribution with rejection threshold, without post-processing. The generator-classes are represented by the three-letter codes defined in Section II-B.

TABLE IV: CRR and average CCR of XceptionNet, DE-FAKE, and VTC in open set attribution by multi-class classification with rejection.

Metric(Step)	XceptionNet	DE-FAKE	ViT
Avg. CCR (Plain)	1.00	0.86	0.93
Avg. CCR (1 Step)	0.79	0.77	0.91
Avg. CCR (2 Steps)	0.61	0.72	0.89
Avg. CCR (3 Steps)	0.43	0.65	0.87
CRR (Plain)	0.70	0.52	0.37
CRR (1 Step)	0.67	0.52	0.26
CRR (2 Steps)	0.66	0.49	0.20
CRR (3 Steps)	0.64	0.49	0.17

IV. CONCLUSIONS

We introduced WILD: a new in-the-wild Image Linkage Dataset for synthetic image attribution. The dataset consists of a closed set and an open set, composed of 10,000 synthetic images each, enabling the evaluation of synthetic source attribution models in both closed-set and open-set scenarios. Furthermore, 5,000 images from each set were post-processed with one, two, and three operators, simulating real-world transformations that images may undergo. To assess WILD, we evaluated seven common source attribution methodologies, providing a benchmark for future research on this dataset. Our findings show that, while synthetic image source attribution is a relatively easy task in a closed set without post-processing, the problem is far from being solved in real-world scenarios. In-the-wild conditions, where unknown generators and

post-processing steps introduce significant variability, led to substantial performance degradation. We believe that WILD represents a challenging dataset for benchmarking the next generation of image source attribution models: It allows to develop models that can work in the wild and tackle real-world attribution tasks. Moreover, as underlined by robustness evaluations, model robustness can be tested and developed on the post-processed images, which proved to be a challenge for most baselines. WILD also opens rooms for interesting future directions: for example, it allows to investigate how prompt semantics affect different generators, but also to conduct explainability studies on source attribution methods.

ACKNOWLEDGEMENT

This work was partially supported by project SERICS (PE00000014) under the MUR National Recovery and Resilience Plan funded by the European Union - NextGenerationEU, and by project FOSTERER, funded by MUR within the PRIN 2022 program under contract 202289RHP.

REFERENCES

- [1] M. Casu, L. Guarnera, P. Caponnetto, and S. Battiato, “GenAI mirage: The impostor bias and the deepfake detection challenge in the era of artificial illusions,” *Forensic Science International: Digital Investigation*, vol. 50, p. 301795, 2024.
- [2] R. Corvi, D. Cozzolino, G. Poggi, K. Nagano, and L. Verdoliva, “Intriguing properties of synthetic images: from generative adversarial networks to diffusion models,” *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 973–982, 2023.

- [3] A. U. Hirte, M. Platscher, T. Joyce, J. Heit, E. Tranvinh, and C. Federau, “Realistic generation of diffusion-weighted magnetic resonance brain images with deep generative models.” *Magnetic resonance imaging*, 2021.
- [4] B. Khoo, R. C.-W. Phan, and C.-H. Lim, “Deepfake attribution: On the source identification of artificially generated images,” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 12, no. 3, p. e1438, 2022.
- [5] L. Bindini, G. Bertazzini, D. Baracchi, D. Shullani, P. Frasconi, and A. Piva, “Tiny autoencoders are effective few-shot generative model detectors,” in *2024 IEEE International Workshop on Information Forensics and Security (WIFS)*. IEEE, 2024, pp. 1–6.
- [6] D. Beniaquev, “Synthetic Faces High Quality - Text 2 Image (SFHQ-T2I) Dataset,” 2024. [Online]. Available: <https://github.com/SelfishGene/SFHQ-T2I-dataset>
- [7] M. A. Rahman, B. Paul, N. H. Sarker, Z. I. A. Hakim, and S. A. Fattah, “Artifact: A large-scale dataset with artificial and factual images for generalizable and robust synthetic image detection,” in *2023 IEEE International Conference on Image Processing (ICIP)*, 2023, pp. 2200–2204.
- [8] M. Zhu, H. Chen, Q. Yan, X. Huang, G. Lin, W. Li, Z. Tu, H. Hu, J. Hu, and Y. Wang, “Genimage: A million-scale benchmark for detecting ai-generated image,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 77771–77782, 2023.
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [10] Adobe, *Adobe Firefly*, 2023, <https://firefly.adobe.com/>.
- [11] OpenAI, “Improving image generation with better captions,” *OpenAI Technical Report*, 2024, accessed: 2025-03-20. [Online]. Available: <https://cdn.openai.com/papers/dall-e-3.pdf>
- [12] B. F. Labs, “Flux,” <https://github.com/black-forest-labs/flux>, 2024.
- [13] Black Forest Labs, “FLUX 1.1 [pro]: Advanced Text-to-Image Generation Model,” 2024, accessed: 2025-03-20. [Online]. Available: <https://blackforestlabs.ai/1-1-pro/>
- [14] Freepik, “Freepik AI Image Generator,” 2024, accessed: 2025-03-20. [Online]. Available: <https://docs.freepik.com/api-reference/mystic/post-mystic>
- [15] Leonardo AI, “Leonardo AI: AI-Powered Creative Image Generation Platform,” 2024, accessed: 2025-03-20. [Online]. Available: <https://leonardo.ai>
- [16] MidJourney, “MidJourney: An AI-powered image generation tool,” 2024, accessed: 2025-03-20. [Online]. Available: <https://www.midjourney.com>
- [17] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [18] S. AI, “Stable diffusion 3.5-large,” <https://huggingface.co/stabilityai/stable-diffusion-3.5-large>, accessed: February 28, 2025.
- [19] H. Face, “Using diffusers: Sdxl turbo,” 2023, accessed: 2025-03-11. [Online]. Available: https://huggingface.co/docs/diffusers/en/using-diffusers/sdxl_turbo
- [20] S. AI, “Starry ai,” 2023, accessed: 2025-03-11. [Online]. Available: <https://starryai.com/>
- [21] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, “Zero-shot text-to-image generation,” *arXiv preprint arXiv:2102.12092*, 2021. [Online]. Available: <https://arxiv.org/abs/2102.12092>
- [22] DeepAI, “DeepAI Text-to-Image Generator,” 2024, accessed: 2025-03-20. [Online]. Available: <https://deepai.org/machine-learning-model/text2img>
- [23] HotPot AI, “HotPot AI: AI-Powered Image and Text Generation Tools,” 2024, accessed: 2025-03-20. [Online]. Available: <https://hotpot.ai/>
- [24] E. Xie, J. Chen, J. Chen, H. Cai, H. Tang, Y. Lin, Z. Zhang, M. Li, L. Zhu, Y. Lu, and S. Han, “SANA: Efficient high-resolution image synthesis with linear diffusion transformers,” *arXiv preprint arXiv:2410.10629*, 2024, accessed: 2025-03-20. [Online]. Available: <https://arxiv.org/abs/2410.10629>
- [25] P. Pernias, D. Rampas, M. L. Richter, C. J. Pal, and M. Aubreville, “Wuerstchen: An efficient architecture for large-scale text-to-image diffusion models,” 2023.
- [26] H. Chefer, Y. Alaluf, Y. Vinker, L. Wolf, and D. Cohen-Or, “Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models,” 2023. [Online]. Available: <https://arxiv.org/abs/2301.13826>
- [27] T. Karras, S. Laine, and T. Aila, “A style-based generator architecture for generative adversarial networks,” 2019. [Online]. Available: <https://arxiv.org/abs/1812.04948>
- [28] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, “Analyzing and improving the image quality of stylegan,” 2020. [Online]. Available: <https://arxiv.org/abs/1912.04958>
- [29] T. Karras, M. Aittala, S. Laine, E. Härkönen, J. Hellsten, J. Lehtinen, and T. Aila, “Alias-free generative adversarial networks,” *Advances in neural information processing systems*, vol. 34, pp. 852–863, 2021.
- [30] Z. Li, J. Zhang, Q. Lin, and et al., “Hunyuan-dit: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding,” 2024.
- [31] G. Gouvine, “Ninasr: Efficient small and large convnets for super-resolution,” <https://github.com/Coloquinte/torchSR/blob/main/doc/NinaSR.md>, 2021.
- [32] N. Ahn, B. Kang, and K.-A. Sohn, “Fast, accurate, and lightweight super-resolution with cascading residual network,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 252–268.
- [33] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, “Enhanced deep residual networks for single image super-resolution,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 136–144.
- [34] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, “Residual dense network for image super-resolution,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2472–2481.
- [35] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, “Image super-resolution using very deep residual channel attention networks,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 286–301.
- [36] J. Ascenso, E. Alshina, and T. Ebrahimi, “The JPEG AI standard: Providing efficient human and machine visual data consumption,” *IEEE Multimedia*, vol. 30, no. 1, pp. 100–111, 2023.
- [37] 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 *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [38] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [39] 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 *Proceedings of the 2023 ACM SIGSAC conference on computer and communications security*, 2023, pp. 3418–3432.
- [40] J. Li, D. Li, C. Xiong, and S. Hoi, “Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation,” in *International conference on machine learning*. PMLR, 2022, pp. 12 888–12 900.
- [41] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.
- [42] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
- [43] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [44] S. Mandelli, P. Bestagini, and S. Tubaro, “When synthetic traces hide real content: Analysis of stable diffusion image laundering,” in *2024 IEEE International Workshop on Information Forensics and Security (WIFS)*. IEEE, 2024, pp. 1–6.
- [45] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly et al., “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [46] T. Ridnik, E. Ben-Baruch, A. Noy, and L. Zelnik-Manor, “Imagenet-21k pretraining for the masses,” 2021. [Online]. Available: <https://arxiv.org/abs/2104.10972>
- [47] D. Cozzolino, G. Poggi, R. Cory, 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.