

Model-agnostic Origin Attribution of Generated Images with Few-shot Examples

Fengyuan Liu¹, Haochen Luo¹, Yiming Li², Philip Torr¹, and Jindong Gu^{1*}

¹ University of Oxford, Oxford OX1 3PJ, UK

² Nanyang Technological University, Singapore 639798, Singapore
oxfengyuan@gmail.com, jindong.gu@outlook.com

Abstract. Recent progress in visual generative models enables the generation of high-quality images. To prevent the misuse of generated images, it is important to identify the origin model that generates them. In this work, we study the origin attribution of generated images in a practical setting where only a few images generated by a source model are available and the source model cannot be accessed. The goal is to check if a given image is generated by the source model. We first formulate this problem as a few-shot one-class classification task. To solve the task, we propose OCC-CLIP, a CLIP-based framework for few-shot one-class classification, enabling the identification of an image’s source model, even among multiple candidates. Extensive experiments corresponding to various generative models verify the effectiveness of our OCC-CLIP framework. Furthermore, an experiment based on the recently released DALL-E-3 API verifies the real-world applicability of our solution.

Keywords: Model Attribution · Generated Images · CLIP Classification

1 Introduction

Recent visual generative models are capable of producing images of exceptional quality, which have raised public concerns regarding Intellectual Property (IP) protection and accountability for misuse. In response to both the challenges and opportunities posed by AIGC, a recent U.S. executive order [3] mandates that all AI-generated content must clearly label its source, such as Stable Diffusion [35]. This makes the attribution of origins for generated images crucial in real-world applications, referring to the process of identifying whether a given image is generated by a particular model.

To address the origin attribution problem above, three main methods have been explored in the community. The first method involves watermarking [26, 32, 46, 48], which requires additional modifications to the generated results, affecting the quality of generation. The second method involves injecting fingerprints [8, 56–58] into the model during training and employing a supervised classifier to identify these fingerprints. This process necessitates changes in training.

* Corresponding author

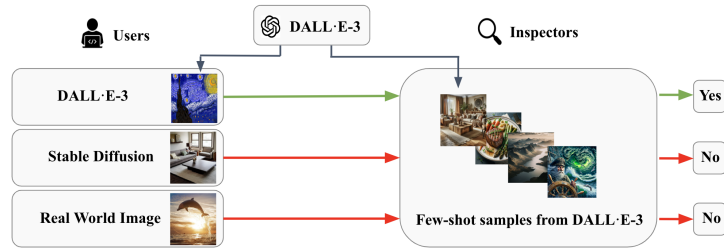


Fig. 1: A simple demonstration of origin attribution in a practical, open-world setting. Inspectors receive a few samples from DALL-E-3. Users then submit an image, which could either be a real photograph or generated by DALL-E-3 or other models. If it’s determined that the image and the provided samples were generated by the same model, we can then identify DALL-E-3 as the origin model of the query image.

Alternation-free approaches have also been proposed, which do not require alterations to the generation or training processes. Specifically, the existing methods utilize inverse engineering [22, 53], based on the idea that a synthetic sample can be most accurately reconstructed by the generator that created it. However, inverse engineering approaches necessitate access to the target model and require sampling many images as references.

In this work, we aim to conduct origin attribution in a practical open-world setting (Fig. 1), where model parameters cannot be accessed and only a few samples generated by the model are available. This setting is meaningful in real-world applications since current generative models, e.g., DALL-E-3 [2], are not open-sourced, and sampling many images from them requires substantial costs.

To overcome the challenges in this setting, we first formulate the problem as a few-shot one-class classification task. Then, we propose a CLIP-based framework as an effective solution, dubbed OCC-CLIP. With our framework, we can determine if a given image and the few-shot available images were generated from the same model. If so, we can confidently identify the model that generated the few images as the origin model of the given image. Furthermore, we also demonstrate that our method can be extended to conduct origin attribution for multiple source models via One-vs-Rest.

Extensive experiments on various generative models have verified that our proposed framework can effectively determine the origin attribution of a given image. Additionally, our framework demonstrates superiority when different numbers of shots are available, and when image preprocessing is applied to the given images. It also proves effective in multi-source origin attribution scenarios. Furthermore, our experiments, based on the recently released DALL-E-3 [2] API, confirm the effectiveness of our solution in real-world commercial systems.

Our main contributions can be summarized as follows.

- We propose a new task within a practical setting where generated images are attributed to the origin model only with few-shot available images generated by the model.
- We formulate the problem as a few-shot one-class classification task and then propose a CLIP-based framework, named OCC-CLIP, to address it.

- Extensive experiments are conducted on 8 generative models, including diffusion models and GANs. Further verification of our solution is carried out on a real-world image generative system, namely, DALL·E-3 [2].

2 Related Work

Deep Visual Generative Models: The advent of deep learning has brought significant advancements in deep visual generative models, leading to the development of sophisticated methods for creating synthetic media. Variational Autoencoders (VAEs) [21, 30, 34] are notable for their dual-structure framework. In VAEs, an encoder condenses complex data into simpler latent representations, which a decoder then uses to reconstruct the original data. Generative Adversarial Networks (GANs) [1, 11, 17, 38, 49] operate on a principle of competition between two components: a generator that creates data samples and a discriminator that judges their authenticity. Through iterative training, both components enhance their capabilities, improving the overall quality of the generated data. Diffusion Models [13, 15, 29, 35, 37, 45] employ a two-phase process. Initially, they transform data into noise, and then methodically remove this noise to reverse the process. This approach can generate highly realistic images.

Origin Attribution of Generated images: Origin attribution, distinct from generated image detection, seeks to determine whether specific images were produced by a particular model. Various strategies have been proposed to address this challenge. One method involves embedding watermarks [26, 32, 46, 48] in images to trace their origins. However, this approach faces limitations as different models might use identical watermarks, which can also be manipulated or removed. Alternatively, injecting unique fingerprints [8, 56–58] into a model during its training phase enables the identification of these markers using a supervised classifier. Although effective, this technique necessitates modifications to the training process and the model’s architecture. A different, modification-free method is inverse engineering [22, 53], which leverages the principle that a synthetic sample can be most accurately reconstructed by its original generator. However, this approach requires access to the target model and extensive sampling of images for comparison. Furthermore, two studies have explored this area under different settings. One develops a multi-class classifier known as an attributor for fake image attribution. Yet, it is only applicable to text-to-image models and operates within a limited, closed-world scenario [43]. Another study [20] demonstrates a theoretical lower bound on attribution accuracy using smaller datasets (e.g. MNIST [23]). This method is only suitable for unconditional GANs and reveals limited scalability to more complex systems, such as DALL·E-3 [2].

Few-shot One-class classification: One-class classification has long been recognized as a challenging problem, with numerous studies [6, 36, 39, 40] addressing it. However, few works have explored the few-shot one-class classification (FS-OCC) problem in the image domain. Previous FS-OCC research in the image domain, specifically based on meta-learning [10], suffers from memory inefficiency, making it unsuitable for processing high-resolution images. In contrast,

the CLIP-based classifier, CoOp [59], has demonstrated excellent performance in few-shot learning settings, where a few images from each class are available. Our setting, however, differs as it involves only a few images from a single class, generated by a generative model. Therefore, instead of pursuing multi-class classification, we utilize CLIP for one-class classification.

3 Approach

In this section, we first introduce the standard CLIP-based classification framework and then present our OCC-CLIP framework for few-shot one-class classification. At the end, we show how to extend our framework to multiple classes.

3.1 Background of CLIP-based Classification

CLIP is a multimodal model pre-trained to predict whether an image matches a text prompt. It includes an image encoder $E_v(\cdot)$ and text encoder $E_t(\cdot)$. The pre-trained CLIP can perform zero-shot multi-class classification by comparing the image with a list of prompts, each representing a class. Formally, assume we have K classes. Let $X^v \in \mathcal{X}$ denote an image and X_i^t represent the i _{th} prompt representing the i _{th} class. The predicted i _{th} class probability of the image X^v is computed as follows:

$$p(\text{class} = i|v) = \frac{\exp(\text{sim}(E_t(X_i^t), E_v(X^v)))}{\sum_{j=1}^K \exp(\text{sim}(E_t(X_j^t), E_v(X^v)))}, \quad (1)$$

where $\text{sim}(\cdot)$ measures the distance between two embeddings, e.g., dot product.

In the classification above, hand-crafted text prompts are applied to represent classes. The prompt designs require specialized knowledge and are time-consuming to create. To alleviate this, Zhou et al. [59] introduced the concept of Context Optimization (CoOp), which employs learnable vectors to refine prompt-related words. Instead of manually designing a prompt, CoOp enables the model to automatically optimize for a suitable prompt. Formally, the prompt for the i _{th} class can be represented as $X_i^p = [t] \otimes [CLASS_i]$, where t is the learnable context vectors, \otimes is a concatenation operation and $CLASS_i$ corresponds to the name of i _{th} class. The objective of optimization is to minimize the error of predicting ground truth Y_i for each X_i^v . This is achieved by using a cross-entropy loss function \mathcal{L} with respect to the learnable prompts. Let f represent the pre-trained CLIP model. The optimization can be described as:

$$\min_{X^p} \sum_{j=1}^N \sum_{i=1}^K \mathcal{L}(f(E_v(X_j^v), E_t(X_i^p)), Y_i), \quad (2)$$

where N is the number of images. In practice, CoOp [59] shows that only a small number of labeled examples are required to learn effective prompts. In summary, prompt learning enables CLIP-based Classification effective in the few-shot learning setting.

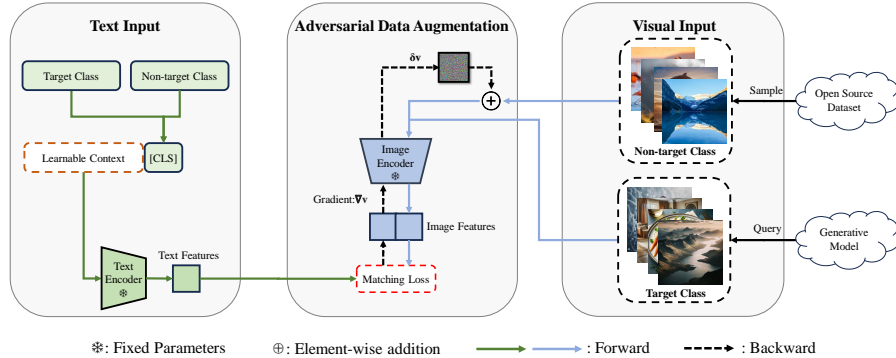


Fig. 2: Overview of OCC-CLIP. The input text is represented by learnable context vectors, followed by two discrete classes: the *target* class corresponds to an image set queried from a generative model, and the *non-target* class corresponds to randomly sourced open-domain images. These classes can be labeled as contrasting pairs, such as non-target vs. target or negative vs. positive. The parameters for the text and image encoders, derived from the CLIP model, are fixed. Adversarial Data Augmentation (ADA) calculates the gradient of each pixel across non-target images. In the training phase, these gradients δ^v are applied to the non-target images.

3.2 CLIP-based Few-shot One-Class Classification

The CLIP-based classifier, CoOp, can achieve excellent performance in the few-shot learning setting where a few images from each class are available. Nevertheless, in our setting, only a few images from one class, generated by a generative model, are available. Hence, a standard CLIP-based classifier cannot be applied directly to solve the few-shot one-class classification task.

We now present our CLIP-based framework for One-Class Classification, called **OCC-CLIP**. In our framework (Fig. 2), the few images collected from a generative model are treated as the target class, while the images randomly sampled from a clean dataset are labeled as the non-target class. The two classes can be labeled as any contrasting pairs, such as non-target vs. target or negative vs. positive. Two learnable prompts corresponding to the target and non-target classes are optimized on these images, respectively.

The few images selected for the non-target class cannot represent the whole distribution of the non-target class well since they are randomly sampled from an open-source dataset (e.g. ImageNet [7]). To overcome this challenge, we propose an adversarial data augmentation (ADA) technique that, during training, extends coverage of the non-target class space and more closely approximates the boundary to the target space, thereby improving the model’s ability to learn the attribution of the target model. ADA aims to maximize the loss by adding small perturbations δ^v to non-target images, while the learnable prompts aim to minimize the loss by learning the boundary between the target and non-target

Algorithm 1: OCC-CLIP Framework

input : X^v : images, X^p : learnable prompts, \mathcal{G} : ground truth, ϵ : updating step size

for *Iterations* **do**

$EV = E_v(X^v), ET = E_t(X^p) ;$	▷ Extract embeddings
$G = \nabla_{\mathcal{X}} \mathcal{L}(f(EV, ET), \mathcal{G}) ;$	▷ Gradient calculation
$\delta^v = \epsilon \cdot \text{sign}(G) ;$	▷ Perturbation generation
$EV \leftarrow E_v(X^v + \delta^v) ;$	▷ Update image embeddings
$X^p \leftarrow \min_{X^p} \mathcal{L}(f(EV, ET), \mathcal{G}) ;$	▷ Prompt optimization

classes. In summary, the optimization in OCC-CLIP can be formulated as:

$$\min_{X^p} \max_{\delta^v} \sum_{j=1}^N \sum_{i=1}^K \mathcal{L}(f(E_v(X_j^v + \delta_j^v), E_t(X_i^p)), Y_i), \quad (3)$$

where δ^v is the adversarial image perturbation computed by our ADA technique.

The implementation of our OCC-CLIP is shown in Algorithm 1. As shown in the algorithm, a forward pass and a backward pass are conducted to obtain gradient information for the images of the non-target class. The gradient information is used to compute adversarial perturbation. The learnable prompt will be updated in another forward and backward passes on the perturbed images. The sensitivity of training hyperparameters is discussed in experiments section.

In the optimization process, both visual and textual encoders of CLIP are frozen. In the verification process, if it is classified into the target class, an image will be determined to be generated by the same generative model as the source model of the target images.

3.3 CLIP-based Few-shot Multi-Class Classification

We also explore the multi-source origin attribution scenarios. For example, to determine if the origin of an image can be attributed to ProGAN [18], Stable Diffusion [35], or Vector Quantized Diffusion [13], we can employ three one-class classifiers corresponding to these models for classification. Given a set of trained K one-class classifiers $\{OCC_1, OCC_2, \dots, OCC_K\}$ for K classes and a threshold θ (e.g. 0.5), for an input sample X^v , let $s_i(X^v)$ denote the score of X^v given by i -th classifier OCC_i . The predicted class $C(X^v)$ for the sample X^v is determined as follows:

$$C(X^v) = \begin{cases} \arg \max_{i \in \{1, \dots, K\}} s_i(X^v) & \text{if } \max_{i \in \{1, \dots, N\}} s_i(X^v) > \theta, \\ \text{others} & \text{otherwise.} \end{cases} \quad (4)$$

Given an image, if the maximum score of X^v , provided by the i -th classifier, exceeds the threshold, then X^v is classified into the i -th class. Otherwise, it is considered to belong to a category outside those defined by the K classifiers.

4 Experiments

In this section, we first describe experimental settings and present our comparison with baseline methods. We also study the sensitivity of our method to various factors, such as target class corresponding to source models, non-target class datasets, the number of available images, and image preprocessing. Furthermore, we show the effectiveness of our framework in multi-source origin attribution scenarios and real-world commercial generation API.

4.1 Experimental Setting

Dataset. There are total 202,520 images generated by five different generative models, i.e. Stable Diffusion Model [35], Latent Diffusion Model [35], GLIDE [29], Vector Quantized Diffusion [13], and GALIP [49], based on the validation set of Microsoft Common Objects in Context (COCO) 2014 dataset [24]. These models are pre-trained on four different datasets, i.e. LAION-5B [41], COCO [24], LAION-400M [42], and filtered CC12M [5]. In total, five image datasets from different source models are generated, namely SD, VQ-D, LDM, Glide, GALIP. To balance the number of image datasets generated by Diffusion models and the number of image datasets generated by GANs, we utilized pre-existing datasets (namely GauGAN [31], ProGAN [18], and StyleGAN2 [19]) as provided by [52]. These datasets collectively serve as a robust benchmark covering two primary generative techniques: GANs and Diffusion Models.

Model. OCC-CLIP utilizes 16 context vectors. This model is built upon the open-source CLIP framework. The image encoder uses the ViT-B/16 architecture. Except for the prompt learner, all pre-trained parameters are fixed. Initial context vectors are stochastically sampled from a Normal distribution characterized by a mean of 0 and a standard deviation of 0.02.

Training Setting. All the generated images are resized to 224×224 and then normalized according to the pre-trained datasets of each model. Stochastic Gradient Descent is utilized as the optimization strategy with a learning rate of 0.0001, modulated through cosine annealing. The cross-entropy loss is utilized as the loss function. By default, the training process is capped at a maximum of 200 epochs for 50-shot scenarios. The test dataset consists of 1,000 images that are randomly selected from the test sets to ensure a reliable evaluation. To counteract the onset of explosive gradients in the nascent phases of training, the learning rate is steadfastly maintained at 1×10^{-5} during the first epoch. The eight generative models (i.e. SD, VQ-D, LDM, Glide, GALIP, ProGAN, StyleGAN2, GauGAN) are iteratively treated as the target class, while four open-source datasets (i.e. COCO [24], ImageNet [7], Flickr [54], and CC12M [5]) are iteratively treated as the non-target class. However, in the default settings, only half of the non-target images are augmented by ADA, SD is designated as the target image set, and COCO is chosen as the non-target image set.

Evaluation. Each model is evaluated by the Area Under the Receiver Operating Characteristic Curve (AUC). The Receiver Operating Characteristic Curve

Table 1: Compare the performance of OCC-CLIP in origin attribution against 12 basic methods. During the training phase, the target class is sourced from SD, and the non-target class is from COCO. In the test phase, the target class remains sourced from SD, but the non-target class is sourced from another generative image dataset. The optimal outcomes for individual datasets are emphasized using **bold** formatting.

Methods	VQ-D	LDM	Glide	GALIP	ProGAN	StyleGAN2	GauGAN	Overall
VGG16 [44]	0.6458	0.5438	0.5652	0.7017	0.6609	0.6170	0.7096	0.6349
ResNet50 [14]	0.6693	0.5485	0.6103	0.7307	0.6646	0.6458	0.7475	0.6595
Inception-v3 [47]	0.6378	0.5349	0.5252	0.6759	0.6453	0.5947	0.6982	0.6160
DenseNet-121 [16]	0.7264	0.5337	0.5730	0.7478	0.7204	0.6645	0.8091	0.6821
ViT-B-16 [9]	0.7502	0.5387	0.6069	0.7837	0.6062	0.6593	0.7259	0.6673
DeiT-B-16 [50]	0.6825	0.5366	0.5698	0.6592	0.5721	0.5898	0.6507	0.6087
CaiT-S-24 [51]	0.6522	0.5301	0.5550	0.6466	0.5697	0.5804	0.6645	0.5998
Swin-B-4 [25]	0.8634	0.7180	0.7473	0.8478	0.5879	0.7054	0.7822	0.7503
Image-Patch [28]	0.6638	0.5269	0.5534	0.6847	0.6624	0.6706	0.7674	0.6470
Feature-Patch [28]	0.6999	0.6116	0.7385	0.6285	0.7136	0.8546	0.8152	0.7231
CLIP [33]	0.7272	0.7243	0.6078	0.7304	0.5229	0.5393	0.6463	0.6426
CoOp [59]	0.9503	0.8373	0.9266	0.9660	0.8861	0.9533	0.9643	0.9263
OCC-CLIP	0.9703	0.8801	0.9519	0.9798	0.9452	0.9651	0.9910	0.9548

(ROC) is a graphical representation that plots the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold levels. The AUC represents the probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one. To reduce randomness, each model is trained 10 times with different training sets each time. Then, the mean AUC and the corresponding standard deviation are reported over the test set. A higher AUC score signifies better performance. For each table, the corresponding standard deviations are shown in the supplementary. More experimental details, such as different testing tasks, can be found in the supplementary.

4.2 Comparison with Baselines

Baselines. Since there are currently no methods perfectly suited to our setting, we conduct a comprehensive evaluation of OCC-CLIP by assessing 12 benchmark methods from various usage areas (see Table 1).

Traditional Binary Classification Models: Four of these methods are notable CNN architectures that are often applied to Computer Vision problems, especially for image classification tasks, including ResNet [14], Inception [47], DenseNet [16], and VGG [44]. The other four models are adaptations of the transformer paradigm, including ViT [9], Deit [50], CaiT [51], and Swin [25]. For both CNN and Transformer models, all layers are retained in a frozen state except for the last classification layer.

Patch-Driven Methods [28]: Two of the methods, Feature-Patch and Image-Patch, are patch-based techniques that are primarily used for deepfake and rely on the ResNet-50 architecture for feature extraction. For the Image-Patch model, each image is divided into 2x2 patches, with each patch serving as a separate input to the ResNet-50 architecture. In the case of the Feature-Patch model,

Table 2: Evaluation sensitivity of OCC-CLIP to Source Models. The leftmost column represents target datasets. In the training phase, the non-target images are from COCO. In the testing phase, the non-target images are from a different generative image dataset shown in the first row. The optimal outcomes for individual datasets are emphasized using **bold** formatting.

Target ↓	Method	SD	VQ-D	LDM	Glide	GALIP	ProGAN	StyleGAN2	GauGAN	Overall
SD	CoOp	–	0.9503	0.8373	0.9266	0.9660	0.8861	0.9533	0.9643	0.9263
	OCC-CLIP	–	0.9703	0.8801	0.9519	0.9798	0.9452	0.9651	0.9910	0.9548
VQ-D	CoOp	0.9988	–	0.7337	0.7435	0.7558	0.9290	0.9924	0.9352	0.8698
	OCC-CLIP	0.9992	–	0.6924	0.7264	0.7327	0.9931	0.9936	0.9936	0.8758
LDM	CoOp	0.9925	0.6793	–	0.6468	0.6263	0.9565	0.9896	0.9758	0.8381
	OCC-CLIP	0.9957	0.7507	–	0.6847	0.6530	0.9956	0.9940	0.9940	0.8676
Glide	CoOp	0.9985	0.8573	0.8314	–	0.6687	0.9629	0.9916	0.9814	0.8988
	OCC-CLIP	0.9998	0.8958	0.8585	–	0.6834	0.9974	0.9949	0.9997	0.9185
GALIP	CoOp	0.9999	0.9036	0.8441	0.7802	–	0.9982	0.9992	0.9996	0.9321
	OCC-CLIP	0.9999	0.9345	0.8626	0.7779	–	0.9999	0.9993	1.0000	0.9392
ProGAN	CoOp	0.9972	0.9475	0.9544	0.9453	0.9818	–	0.9278	0.7993	0.9362
	OCC-CLIP	0.9961	0.9477	0.9585	0.9320	0.9885	–	0.8471	0.8885	0.9369
StyleGAN2	CoOp	0.9992	0.9856	0.9850	0.9584	0.9849	0.8687	–	0.9543	0.9623
	OCC-CLIP	0.9996	0.9937	0.9851	0.9652	0.9926	0.9649	–	0.9946	0.9851
GauGAN	CoOp	0.9991	0.9388	0.9862	0.9757	0.9901	0.6812	0.9676	–	0.9368
	OCC-CLIP	0.9982	0.9132	0.9745	0.9593	0.9958	0.7752	0.9425	–	0.9370

which shares the same backbone, the last five layers of the standard ResNet-50 architecture are discarded and the remaining structures are employed to extract image features. Subsequently, these features are segmented into 4x4 patches, with each patch being passed through the concluding linear classification layer.

Vision-Language Models: Considering the power of the vision-language model in tackling downstream classification tasks, zero-shot CLIP [33] and CoOp [59] are also evaluated. The zero-shot CLIP [33] approach utilizes custom-crafted prompts (hard prompt), adopting the format "a photo of a [CLASS]" to tackle the tasks of origin attribution. Conversely, CoOp employs prompt tuning (soft prompt) and distinguishes itself by using $\{v_1, v_2, \dots, v_{16}, [\text{CLASS}]\}$ where v_i is an adjustable context vector and [CLASS] is the class token which is deliberately located at the end of the sequence.

Results and Analysis. Table 1 provides a comparative performance analysis of OCC-CLIP and 12 benchmark models. During the training phase, the target class is sourced from SD, and the non-target class is from COCO. In the test phase, the target class remains sourced from SD, but the non-target class is sourced from another generative image dataset (VQ-D, LDM, Glide, GALIP, ProGAN, StyleGAN2, and GauGAN). The goal of this section is to assess each model’s ability to accurately determine whether the origin of a given image can be attributed to Stable Diffusion.

Among the CNN-based and Transformer-based baselines, Swin outshines its counterparts. However, the overall disparity in performance between CNN-based and Transformer-based models is not obvious. The following two methods, Image-Patch and Feature-Patch are both based on ResNet50. Image-Patch performs worse than ResNet50 while Feature-Patch has somewhat but limited improvements compared with ResNet50. This suggests that current methods

Table 3: Evaluation of OCC-CLIP on Different Non-target Image Datasets. During training phase, the non-target class is chosen from One of the four datasets (COCO, CC12M, Flickr, and ImageNet). The target dataset is from SD. The performance is averaged across seven different testing tasks.

Methods	COCO	CC12M	Flickr	ImageNet
CoOp	0.9263	0.8689	0.8788	0.9377
OCC-CLIP	0.9548	0.9320	0.9355	0.9710

Table 4: Evaluation of OCC-CLIP with Various ADA Approaches. ‘Non-Target’ applies ADA to the non-target images; ‘T’ applies it to the target images; ‘Both’ applies it to both the target and non-target images; and ‘T-NT’ treats part of the target images as non-target ones post-ADA. The results are averaged across seven testing tasks.

Methods	None	Non-Target	Both	Target	T-NT
OCC-CLIP	0.9263	0.9548	0.9080	0.8236	0.9514

used to do deepfake detection perform poorly in this few-shot origin attribution scenarios. In addition to the aforementioned classical methods, OCC-CLIP surpasses widely-spread vision-language models such as zero-shot CLIP and CoOp. In detail, OCC-CLIP outperforms zero-shot CLIP by an average of around 31% and CoOp by approximately 2.9%. This progression emphasizes the value of ADA. In a comprehensive assessment, OCC-CLIP emerges as a leading standard, underscoring its exceptional ability in origin attribution.

4.3 Ablation Study

Sensitivity to Source Models. Table 2 shows the performance of CoOp and OCC-CLIP in origin attribution with eight different source models. During training, the target dataset is from the dataset shown in the leftmost column, and the non-target dataset is COCO. During testing, the target remains the same, but the non-target dataset is replaced with one of the other generated datasets shown in the first row. For example, the source model for the second row in the table is Stable Diffusion [35]. OCC-CLIP demonstrates superior performance in most cases and outperforms the baseline on average in the results from testing with seven other datasets. Overall, our proposed model has exhibited superior performance and generalization capabilities in determining if an image is from the same source model as a set of target images, compared to the baseline.

Sensitivity to Selection of Non-target Class. As shown in Table 3, we select non-target images from one of four different clean datasets: COCO, ImageNet, Flickr, or CC12M. The target images are from SD. The average AUC score over 7 different testing tasks is computed. It can be noted that although the choice of the source of non-target images can somewhat affect the performance of origin attribution, our framework consistently outperforms the baseline.

Sensitivity to ADA Settings. We employ four methods for applying ADA: 1) ‘Non-Target’: on half of the non-target images; 2) ‘Target’: on half of the target

Table 5: Evaluation of OCC-CLIP Relative to the Proportion of Augmented Non-Target Images. Five proportions are investigated: 0%, 25%, 50%, 75%, and 100%. The results are averaged across seven testing tasks.

Methods	0%	25%	50%	75%	100%
OCC-CLIP	0.9263	0.9405	0.9548	0.9558	0.9172

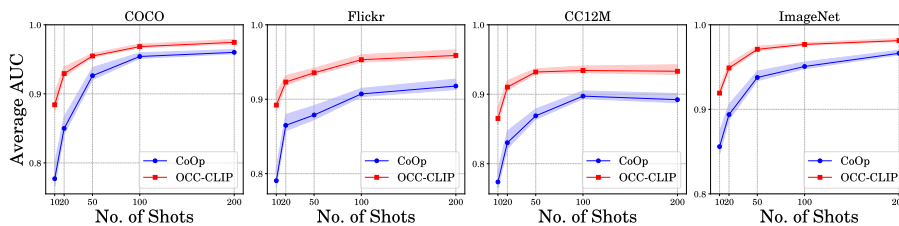


Fig. 3: Evaluation of OCC-CLIP on Different Numbers of Shots. This figure shows the average origin attribution performance of CoOp and OCC-CLIP on 7 different testing tasks with a variable number of shots: 10, 20, 50, 100, and 200. The target dataset is from SD. The non-target datasets are from COCO, CC12M, Flickr, or ImageNet.

images; 3) ‘Both’: on half of both the non-target images and target images; 4) ‘T-NT’: on half of the target images which are then treated as non-target ones. The non-target images are selected from the COCO dataset, while the target images are from SD. As shown in Table 4, it is evident that applying ADA on non-target images yields the best performance. This outcome is reasonable, as gradient ascent on the non-target image set enlarges the learned space of this set and further approximates the boundary to the target image set. Applying ADA on ‘T-NT’ is the next most effective method, performing only slightly worse than ‘Non-Target’. However, the boundary learned by ‘T-NT’ may be too tight for effective origin attribution tasks. Applying ADA on ‘Both’ is less effective than doing nothing, as ‘Both’ also enlarges the learned distribution of the target image set. Solely applying ADA on the target set performs significantly worse than doing nothing. This is because this approach not only enlarges the learned distribution space of the target images but also shrinks the learned distribution space of the non-target image set, leading to a higher likelihood of misclassifying many images as target images.

Sensitivity to the Proportion of Augmented Non-target Images. We further investigate the influence of the proportion of augmented non-target images on the performance of OCC-CLIP. As indicated in Table 5, the performance of origin attribution is optimal when 50% to 75% of non-target images are augmented. However, when all data are augmented, the performance deteriorates, even falling below that with 0% non-target images augmented. This may be because augmenting all data could lead to a loss of the original distribution space of the non-target image set.

Sensitivity to the Number of Shots. We then analyze how the performance of OCC-CLIP relates to the number of shots. Our investigation covers the effects

Table 6: Evaluation of OCC-CLIP under Image Processing. Six image processing methods are executed: Gaussian Blur, Gaussian Noise, Grayscale, Rotation, Flip, and a mixture of all these data augmentation methods. The results are averaged across seven testing tasks.

Method	None	Gaussian Blur	Gaussian Noise	Grayscale	Rotation	Flip	Mixture
CoOp	0.9263	0.8539	0.8859	0.8017	0.9040	0.9241	0.7457
OCC-CLIP	0.9548	0.9002	0.9089	0.8405	0.9337	0.9479	0.7538

Table 7: ADA vs. other data augmentation methods. Different ways of data augmentation are used during training. The results are averaged across seven testing tasks.

Methods	Gaussian Blur	Gaussian Noise	Grayscale	Rotation	Flip	Mixture	None	ADA
OCC-CLIP	0.8639	0.8340	0.8308	0.8465	0.8422	0.7991	0.9263	0.9548

of using 10, 20, 50, 100, and 200 shots. In these experiments, SD is used for the target images, while one of COCO, ImageNet, Flickr, and CC12M is employed as the non-target dataset. Other generated image datasets are used for testing. As shown in Figure 3, OCC-CLIP consistently outperforms CoOp in all scenarios, particularly when the number of shots is small. Additionally, the incremental improvements in AUC tend to diminish as more shots are added, eventually plateauing at an upper boundary when the number of shots reaches 200.

Sensitivity to Image Processing in Verification Stage. Table 6 presents a comparative analysis of the performance between OCC-CLIP and baseline methods in the face of potential image processing: Gaussian Blur, Gaussian Noise, Grayscale, Rotation, Flip, or a mixture of those attacks. The data indicate that OCC-CLIP exhibits superior robustness compared to the baseline across a range of potential processing of input images. This enhanced resilience of OCC-CLIP underscores its effectiveness in safeguarding against various forms of image manipulation, highlighting its utility in origin attribution.

Sensitivity to Choice of Prompts. Since we focus on the one-class classification question, the classes in our setting can be labeled as any contrasting pairs, such as negative vs. positive. Therefore, we explore the effect of the choice of prompts on the performance of OCC-CLIP. As shown in the supplementary material, the choice of prompts can have some different effects. However, regardless of which pair of prompts is chosen, OCC-CLIP always outperforms the baseline.

4.4 Comparison to Standard Data Augmentation

Data augmentation methods have been extensively developed. In real-fake detection tasks, Wang [52] and CR [4] utilized a combination of various data augmentation methods to improve model performance. To conduct a more comprehensive evaluation of ADA, we compare its performance with other common data augmentation methods: Gaussian Blur, Gaussian Noise, Grayscale, Rotation, Flip, and a mixture of all these. According to Table 7, all traditional data augmentation methods perform poorly, and in some cases, even degrade the model’s

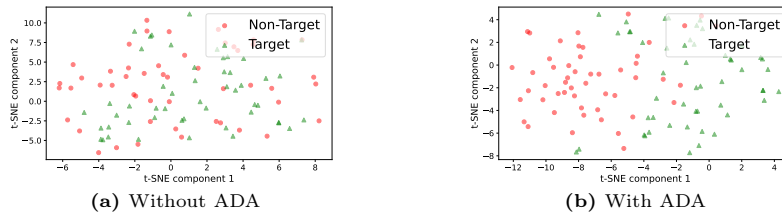


Fig. 4: Visualize the 2-dimensional mapping of image features with t-SNE [27]. The target images are from SD, and the non-target images are from COCO. Figure (a) shows the distribution of non-target images and target images without using ADA. Figure (b) shows the distribution of augmented non-target images and target images after using the ADA technique.

Table 8: Evaluation of OCC-CLIP with commercial generation API. The target images are generated by DALL·E-3. The non-target images are from COCO. The results are averaged across eight testing tasks.

Methods	10	20	30	50
CoOp	0.9216	0.9687	0.9563	0.9788
OCC-CLIP	0.9754	0.9848	0.9936	0.9959

performance. Only ADA shows some improvements. Consequently, we can infer that normal data augmentation methods are not suitable for few-shot one-class classification scenarios.

4.5 Understanding Adversarial Data Augmentation

We utilize the embeddings from CLIP’s image encoder for visualization. According to Figure 4a, it is evident that there is no clear boundary between the target image set and the non-target image set, i.e., images from SD and images from COCO. However, after applying ADA, as illustrated in Figure 4b, a clearer boundary emerges between the augmented non-target image set and the target image set, demonstrating the effectiveness of ADA. During training, unchanged non-target images and augmented non-target images are mixed. This approach not only preserves the properties of the original non-target images but also ensures that the augmented non-target images approximate the boundary of the distribution of the target images.

4.6 Source Model Attribution with Commercial Generation API

To investigate the effectiveness of OCC-CLIP in the real world, we utilize the latest commercial digital image-generating model – DALL·E-3 [2] – to generate images. As shown in the supplementary material, 103 prompts are randomly selected from the annotations of the validation set of COCO [24], resulting in a total of 200 images generated. The default setting of DALL·E-3 is to generate 2 images simultaneously. However, due to OpenAI’s content policy, some prompts

Table 9: Employing an ensemble of one-class classifiers for multi-class classification tasks. The following scenarios are considered: 2 classes, 4 classes, 6 classes, and 8 classes.

Num of Class	Method	LDM	Glide	GALIP	ProGAN	StyleGAN2	GauGAN	Overall
2 classes	CoOp	0.6366	0.6448	0.6378	0.6928	0.8706	0.7259	0.7014
	OCC-CLIP	0.6449	0.6499	0.6457	0.8742	0.8938	0.9255	0.7723
4 classes	CoOp	–	–	0.5863	0.6124	0.7004	0.6376	0.6342
	OCC-CLIP	–	–	0.6130	0.7461	0.7319	0.7908	0.7204
6 classes	CoOp	–	–	–	–	0.6459	0.6092	0.6276
	OCC-CLIP	–	–	–	–	0.6587	0.6727	0.6657
8 classes	CoOp	–	–	–	–	–	–	0.7334
	OCC-CLIP	–	–	–	–	–	–	0.7678

can only generate one image. In this scenario, the targets are images generated by DALL·E-3, while the non-targets are images from other datasets. We compare the performance of detection across different numbers of shots: 10, 20, 30, 50. According to Table 8, our method consistently outperforms the baseline and performs well even with as few as 10 shots.

4.7 Model Attribution with Multiple Source Models

In addition to verifying single models, we also investigate multi-source origin attribution conditions. We explore the process of verifying multiple sources using the eight trained one-class classifiers. Since this involves multi-source origin attribution, the evaluation metric used in this section is **accuracy** instead of AUC. We evaluate the model under the following conditions: 2 classes (SD, VQ-D), 4 classes (SD, VQ-D, LDM, Glide), 6 classes (SD, VQ-D, LDM, Glide, GALIP, ProGAN), and 8 classes (SD, VQ-D, LDM, Glide, GALIP, ProGAN, StyleGAN2, GauGAN). As shown in Table 28, the accuracy decreases with an increasing number of source models, except in the 8 classes condition. We believe this exception for the 8 classes condition occurs because the attribution of each model has been effectively learned by one of the eight classifiers. In every multi-class classification scenario mentioned above, one-class classifiers trained with OCC-CLIP consistently outperform those trained with CoOp.

Furthermore, our methods can be adapted to directly train a multi-class classifier. As demonstrated in the table included in the supplementary material, the multi-class classifier trained using the OCC-CLIP approach is still shown to be more effective compared to the one trained using the CoOp approach.

5 Conclusions

In this work, we study origin attribution in a practical setting where only a few images generated by a source model are available and the source model cannot be accessed. The introduced problem is first formulated as a few-shot one-classification task. A simple yet effective solution CLIP-based framework is proposed to solve the task. Our experiments on both open-source popular

generative models and commercial generation API shows the effectiveness of our framework. Our OCC-CLIP framework can also be applied to solve the few-shot one-classification task in other domains, which we leave in future work. Another future work is to evaluate the adversarial robustness of our framework and build adversarially robust variants.

References

1. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein generative adversarial networks. In: International conference on machine learning. pp. 214–223. PMLR (2017) [3](#)
2. Betker, J., Goh, G., Jing, L., Brooks, T., Wang, J., Li, L., Ouyang, L., Zhuang, J., Lee, J., Guo, Y., et al.: Improving image generation with better captions. Computer Science. <https://cdn.openai.com/papers/dall-e-3.pdf> **2**(3), 8 (2023) [2](#), [3](#), [13](#), [24](#)
3. BIDEN, J.R.: Executive order on the safe, secure, and trustworthy development and use of artificial intelligence (Oct 2023), <https://www.whitehouse.gov/briefing-room/presidential-actions/2023/10/30/executive-order-on-the-safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence/> [1](#)
4. Chandrasegaran, K., Tran, N.T., Binder, A., Cheung, N.M.: Discovering transferable forensic features for cnn-generated images detection. In: European Conference on Computer Vision. pp. 671–689. Springer (2022) [12](#)
5. Changpinyo, S., Sharma, P., Ding, N., Soricut, R.: Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3558–3568 (2021) [7](#), [19](#)
6. Chen, J., Sathe, S., Aggarwal, C., Turaga, D.: Outlier detection with autoencoder ensembles. In: Proceedings of the 2017 SIAM international conference on data mining. pp. 90–98. SIAM (2017) [3](#)
7. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009) [5](#), [7](#)
8. Ding, Y., Thakur, N., Li, B.: Does a gan leave distinct model-specific fingerprints. In: Proceedings of the BMVC (2021) [1](#), [3](#)
9. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020) [8](#), [20](#), [29](#)
10. Frikha, A., Krompaß, D., Köpken, H.G., Tresp, V.: Few-shot one-class classification via meta-learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 7448–7456 (2021) [3](#)
11. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. Communications of the ACM **63**(11), 139–144 (2020) [3](#)
12. Grubbs, F.E.: Sample criteria for testing outlying observations. University of Michigan (1949) [25](#)
13. Gu, S., Chen, D., Bao, J., Wen, F., Zhang, B., Chen, D., Yuan, L., Guo, B.: Vector quantized diffusion model for text-to-image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10696–10706 (2022) [3](#), [6](#), [7](#), [19](#)

14. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016) [8](#), [20](#), [29](#)
15. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. *Advances in neural information processing systems* **33**, 6840–6851 (2020) [3](#)
16. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4700–4708 (2017) [8](#), [20](#), [29](#)
17. Kang, M., Zhu, J.Y., Zhang, R., Park, J., Shechtman, E., Paris, S., Park, T.: Scaling up gans for text-to-image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10124–10134 (2023) [3](#)
18. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196* (2017) [6](#), [7](#), [19](#)
19. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8110–8119 (2020) [7](#), [19](#)
20. Kim, C., Ren, Y., Yang, Y.: Decentralized attribution of generative models. *arXiv preprint arXiv:2010.13974* (2020) [3](#)
21. Kingma, D.P., Welling, M., et al.: An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning* **12**(4), 307–392 (2019) [3](#)
22. Laszkiewicz, M., Ricker, J., Lederer, J., Fischer, A.: Single-model attribution via final-layer inversion. *arXiv preprint arXiv:2306.06210* (2023) [2](#), [3](#), [25](#)
23. LeCun, Y., Cortes, C., Burges, C., et al.: Mnist handwritten digit database (2010) [3](#)
24. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. pp. 740–755. Springer (2014) [7](#), [13](#), [19](#), [24](#)
25. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 10012–10022 (2021) [8](#), [20](#), [29](#)
26. Luo, L., Chen, Z., Chen, M., Zeng, X., Xiong, Z.: Reversible image watermarking using interpolation technique. *IEEE Transactions on information forensics and security* **5**(1), 187–193 (2009) [1](#), [3](#)
27. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. *Journal of machine learning research* **9**(11) (2008) [13](#)
28. Mandelli, S., Bonettini, N., Bestagini, P., Tubaro, S.: Detecting gan-generated images by orthogonal training of multiple cnns. In: 2022 IEEE International Conference on Image Processing (ICIP). pp. 3091–3095. IEEE (2022) [8](#), [29](#)
29. Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M.: Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741* (2021) [3](#), [7](#), [19](#)
30. Oussidi, A., Elhassouny, A.: Deep generative models: Survey. In: 2018 International conference on intelligent systems and computer vision (ISCV). pp. 1–8. IEEE (2018) [3](#)

31. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2337–2346 (2019) 7, 19
32. Pereira, S., Pun, T.: Robust template matching for affine resistant image watermarks. *IEEE transactions on image Processing* **9**(6), 1123–1129 (2000) 1, 3
33. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) 8, 9, 29
34. Rezende, D.J., Mohamed, S., Wierstra, D.: Stochastic backpropagation and approximate inference in deep generative models. In: International conference on machine learning. pp. 1278–1286. PMLR (2014) 3
35. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022) 1, 3, 6, 7, 10, 19
36. Sabokrou, M., Khalooei, M., Fathy, M., Adeli, E.: Adversarially learned one-class classifier for novelty detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3379–3388 (2018) 3
37. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* **35**, 36479–36494 (2022) 3
38. Sauer, A., Karras, T., Laine, S., Geiger, A., Aila, T.: Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. *arXiv preprint arXiv:2301.09515* (2023) 3
39. Schlegl, T., Seeböck, P., Waldstein, S.M., Schmidt-Erfurth, U., Langs, G.: Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In: International conference on information processing in medical imaging. pp. 146–157. Springer (2017) 3
40. Schölkopf, B., Platt, J.C., Shawe-Taylor, J., Smola, A.J., Williamson, R.C.: Estimating the support of a high-dimensional distribution. *Neural computation* **13**(7), 1443–1471 (2001) 3
41. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems* **35**, 25278–25294 (2022) 7, 19
42. Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., Komatsuzaki, A.: Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114* (2021) 7, 19
43. Sha, Z., Li, Z., Yu, N., Zhang, Y.: 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. pp. 3418–3432 (2023) 3
44. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014) 8, 20, 29
45. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502* (2020) 3

46. Swanson, M.D., Zhu, B., Tewfik, A.H.: Transparent robust image watermarking. In: Proceedings of 3rd IEEE International Conference on Image Processing. vol. 3, pp. 211–214. IEEE (1996) [1](#), [3](#)
47. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2818–2826 (2016) [8](#), [20](#), [29](#)
48. Tancik, M., Mildenhall, B., Ng, R.: Stegastamp: Invisible hyperlinks in physical photographs. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2117–2126 (2020) [1](#), [3](#)
49. Tao, M., Bao, B.K., Tang, H., Xu, C.: Galip: Generative adversarial clips for text-to-image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14214–14223 (2023) [3](#), [7](#), [19](#)
50. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: International conference on machine learning. pp. 10347–10357. PMLR (2021) [8](#), [20](#), [29](#)
51. Touvron, H., Cord, M., Sablayrolles, A., Synnaeve, G., Jégou, H.: Going deeper with image transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 32–42 (2021) [8](#), [20](#), [29](#)
52. Wang, S.Y., Wang, O., Zhang, R., Owens, A., Efros, A.A.: Cnn-generated images are surprisingly easy to spot... for now. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8695–8704 (2020) [7](#), [12](#), [19](#)
53. Wang, Z., Chen, C., Zeng, Y., Lyu, L., Ma, S.: Alteration-free and model-agnostic origin attribution of generated images. arXiv preprint arXiv:2305.18439 (2023) [2](#), [3](#), [25](#)
54. Young, P., Lai, A., Hodosh, M., Hockenmaier, J.: From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. Transactions of the Association for Computational Linguistics **2**, 67–78 (2014) [7](#)
55. Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., Xiao, J.: Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. arXiv preprint arXiv:1506.03365 (2015) [19](#)
56. Yu, N., Davis, L.S., Fritz, M.: Attributing fake images to gans: Learning and analyzing gan fingerprints. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 7556–7566 (2019) [1](#), [3](#)
57. Yu, N., Skripniuk, V., Abdelnabi, S., Fritz, M.: Artificial fingerprinting for generative models: Rooting deepfake attribution in training data. In: Proceedings of the IEEE/CVF International conference on computer vision. pp. 14448–14457 (2021) [1](#), [3](#)
58. Yu, N., Skripniuk, V., Chen, D., Davis, L., Fritz, M.: Responsible disclosure of generative models using scalable fingerprinting. arXiv preprint arXiv:2012.08726 (2020) [1](#), [3](#)
59. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. International Journal of Computer Vision **130**(9), 2337–2348 (2022) [4](#), [8](#), [9](#), [29](#)

A Generation of Datasets

There are a total of 202,520 fake images generated by five different generative models, namely, Stable Diffusion Model [35], Latent Diffusion Model [35], GLIDE [29], Vector Quantized Diffusion [13], and GALIP [49]. The prompts used are the first captions of each image in the validation set of the Microsoft Common Objects in Context (COCO) 2014 dataset [24]. These models were pre-trained on four different datasets: LAION-5B [41], COCO [24], LAION-400M [42], and filtered CC12M [5]. In total, as shown in Figure 5, five synthetic image datasets are generated, namely SD, VQ-D, LDM, GLIDE, and GALIP.

SD: This marks the data generated by the Stable Diffusion Model [35]³. It was pre-trained on the LAION-5B dataset. The size of the generated images is 512×512 . The pre-trained model used is ‘sd-v1-4’.

VQ-D: This marks the data generated by the Vector Quantized Diffusion (VQ-D) Model [13]⁴. It was pre-trained on the COCO dataset. The size of the generated images is 256×256 . The pre-trained model ‘coco_pretrained’ served as the backbone of this method.

LDM: This marks the data generated by the Latent Diffusion Model [35]⁵. It was pre-trained on the LAION-400M dataset. The size of the generated images is 256×256 . The pre-trained model, ‘txt2img-f8-large’, is utilized.

GLIDE: This marks the data generated by GLIDE [29]⁶. It was pre-trained on filtered CC12M. The size of the generated images is 256×256 .

GALIP: This marks the dataset generated by GALIP [49]⁷. The model was pre-trained on the COCO dataset. The size of the generated images is 256×256 .

We also utilized pre-existing datasets (namely GauGAN [31], ProGAN [18], and StyleGAN2 [19]) as provided by [52]. GauGAN was trained on the COCO [24] dataset, while ProGan and StyleGan2 were trained on the LSUN [55] dataset. The size of images from those datasets is 256×256 .

B Implementation Details

For all models, we randomly selected 1,000 images from each dataset for testing. Besides, 500 images were chosen from each dataset for training, divided into 10 sets: train1, train2, ..., train10. Each training set contained 50 images from a non-target dataset and 50 from a target dataset. After training for 200 epochs, we evaluated the results on a test set comprising 1,000 testing images from both the non-target and target datasets. There was no dedicated validation set in our experimental setup. By default, adversarial data augmentation was applied to half of the non-target images, with a perturbation step size, denoted as ϵ , of 0.1. The Average AUC value was calculated as the mean of the AUC values obtained

³ <https://github.com/CompVis/stable-diffusion>

⁴ <https://github.com/microsoft/VQ-Diffusion>

⁵ <https://github.com/CompVis/latent-diffusion>

⁶ <https://github.com/openai/glide-text2im>

⁷ <https://github.com/tobran/GALIP>

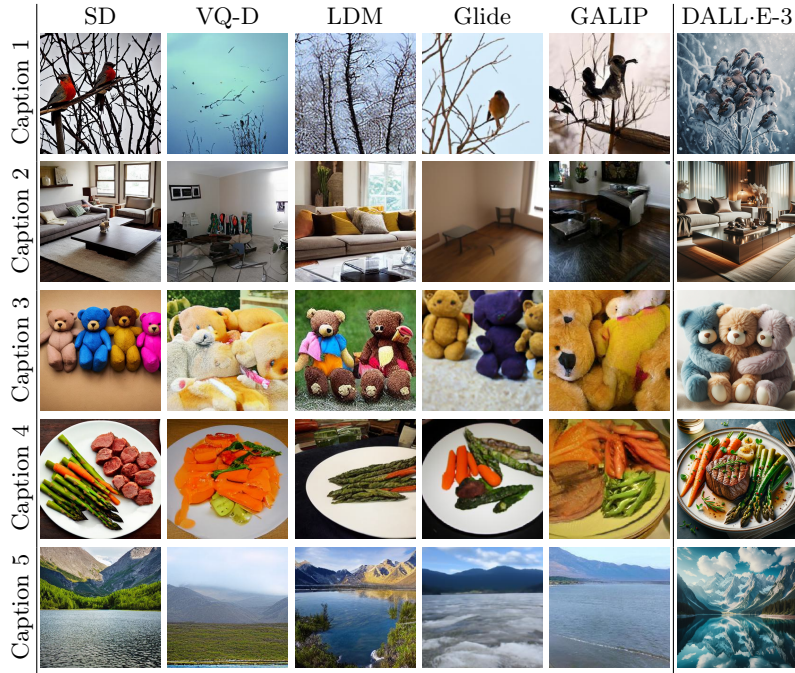


Fig. 5: This figure presents a demo of synthesized images produced by six distinct models, each predicated upon four specific captions. Caption 1: *Birds perch on a bunch of twigs in the winter.* Caption 2: *A coffee table sits in the middle of a living room.* Caption 3: *Three teddy bears, each a different color, snuggling together.* Caption 4: *The dinner plate has asparagus, carrots and some kind of meat.* Caption 5: *A large body of water sitting below a mountain range.*

from the 10 trained models, each trained on a distinct training set and tested on the same test set.

C Comparison with Baselines

The standard deviation of each mean of AUC is shown in Table 10 and Table 11. The reason why the CLIP model has a 0 standard deviation is that it is a zero-shot model, which means there are no parameters changed. Therefore, when testing on the same testing set, the result will not change.

ResNet [14], Inception [47], DenseNet [16], VGG [44], ViT [9], Deit [50], Cait [51], and Swin [25] are either milestones or leading-edge models in the field of Computer Vision. For the Image-Patch model, each image is subdivided into 2x2 patches, with each patch serving as a separate input to the ResNet-50 architecture. The ultimate prediction from the model is the average of the predictions made for each of these patches. In the case of the Feature-Patch model, which shares the same backbone, the last five layers of the standard

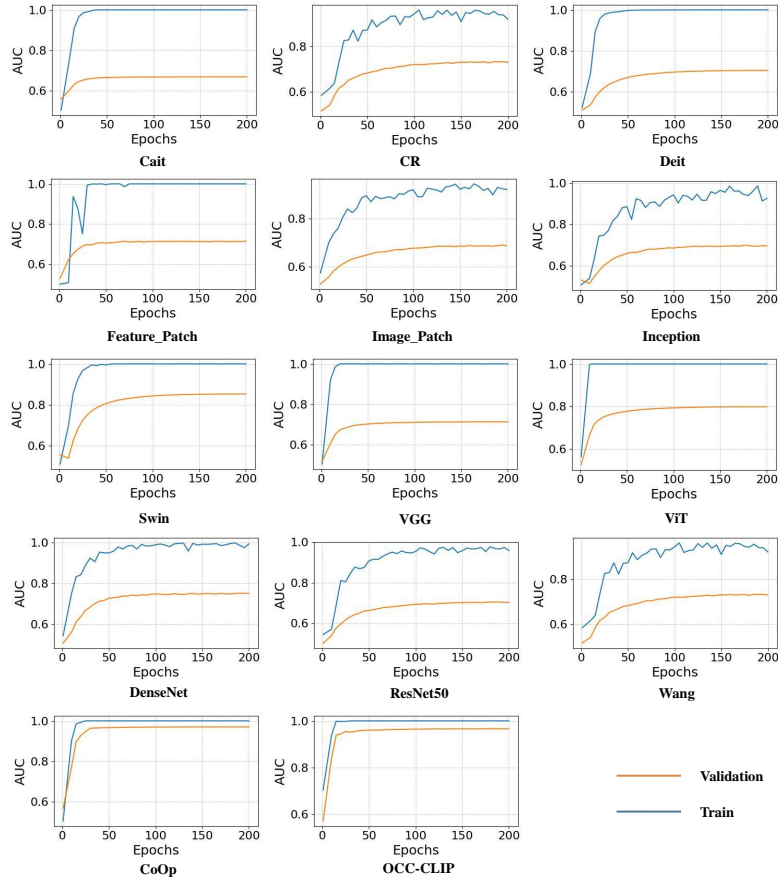


Fig. 6: This figure illustrates the training and validation Area Under the Curve (AUC) metrics for 14 models (13 baselines + OCC-CLIP), comparing their performance in terms of both training accuracy and validation reliability.

ResNet-50 architecture are discarded, and the remaining structure is employed to extract image features. Subsequently, these features are segmented into 4x4 patches, with each patch being passed through the final linear classification layer. The model’s final output is the average prediction over these patches.

To visualize the training loss, training AUC, validation loss, and validation AUC, apart from images in the training set and test set, I randomly select 1000 images from both non-target and target datasets to build a validation set. Figure 6 shows the change of loss of the train set and validation set with the increasing number of epochs. Figure 7 shows the change of AUC of the train set and validation set with the increasing number of epochs. The reason why the loss of the train set is larger than the loss of the validation set in the Image_Patch model is that during training, each image is subdivided into 4 patches.

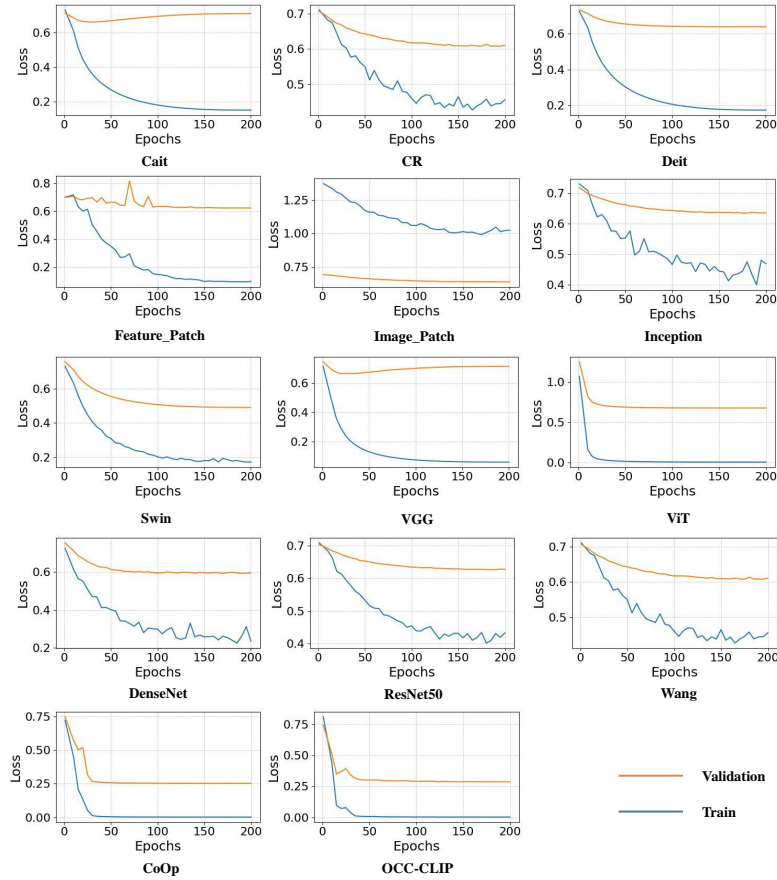


Fig. 7: This figure illustrates the comparative analysis of training and validation loss across 14 models (13 baselines + OCC-CLIP). It provides a visual representation of how each model’s loss metrics evolve over the course of training.

D Sensitivity to Source Models and Value of Epsilon

To test the effect of the choice of ϵ on source models, we conducted a comprehensive evaluation using eight models with six different values of ϵ . According to Table 12 and Table 13, it can be concluded that varying the step size ϵ has different effects on source model attribution. With the adjustment of ϵ , it was observed that the use of adversarial data augmentation can have a positive effect on source model attribution tests.

E Sensitivity to Selection of Non-target Class and ADA Settings

Table 14 and Table 15 show the sensitivity evaluation of various ADA methods on different open-world real image datasets. In the default setting, the target dataset is SD, and the non-target is COCO. ADA is applied to half of the non-target images. 1) During the training phase, conditioned on applying ADA to only half of the non-target images, three additional non-target image datasets are used: ImageNet, Flickr, and CC12M. 2) When selecting non-target images exclusively from COCO, ADA is applied to half of the target images, to half of both non-target and target images, and to half of the target images, which are then treated as non-target. It is observed that the choice of non-target images affects the performance of OCC-CLIP. Additionally, applying ADA to non-target images is most effective.

F Sensitivity to the Proportion of Augmented Non-target Images

Table 16 and Table 17 present the standard deviation when applying ADA to varying proportions of non-target images. Results from seven different testing tasks and their average are also included.

G Sensitivity to the Number of Shots

Table 18 and Table 19 illustrate the mean origin attribution performance of CoOp and OCC-CLIP across seven different testing tasks with varying numbers of shots: 10, 20, 30, 40, 50, 100, and 200. In the training phase, SD is used as the target dataset, and one of COCO, CC12M, Flickr, or ImageNet is used as the non-target dataset.

H Sensitivity to Image Processing in Verification Stage

Tables 20 and 21 present a comparative analysis of the performance between OCC-CLIP and baseline methods in response to potential image processing techniques: Gaussian Blur, Gaussian Noise, Grayscale, Rotation, Flip, or a combination of these attacks. The standard deviations and the testing results across seven different tasks, along with their average, are shown. The data indicate that OCC-CLIP exhibits superior robustness compared to the baseline in most cases, demonstrating enhanced resilience to various forms of image manipulation. This underscores the effectiveness of OCC-CLIP in safeguarding against image-based copyright infringements.

I Sensitivity to Choice of Prompts

We explore the impact of prompt selection on the performance of OCC-CLIP. As indicated in Tables 22 and 23, the choice of prompts influences performance across seven different testing tasks. However, OCC-CLIP consistently outperforms the baseline regardless of the prompt pair chosen.

J Comparison to Standard Data Augmentation

In a comprehensive evaluation of ADA, we compare its effectiveness with common data augmentation methods: Gaussian Blur, Gaussian Noise, Grayscale, Rotation, Flip, and combinations of these. According to Tables 24 and 25, traditional data augmentation methods generally perform poorly, often degrading performance across many of the seven testing tasks. Only ADA shows improvements, suggesting that standard data augmentation methods may not be suitable for few-shot one-class classification scenarios.

K Source Model Attribution with Commercial Generation API

To assess OCC-CLIP’s real-world effectiveness, we used the latest commercial digital image-generating model, DALL·E-3 [2], to generate images. As detailed in Section N, 103 prompts were randomly selected from the annotations of COCO’s validation set [24]. In addition to the images generated from these prompts, additional sample images are presented in Figure 5. During the training phase, the target images are those generated by DALL·E-3, while the non-target images are from COCO. In the testing phase, the non-target images are from one of the other generated datasets. The performance of detection across various shot numbers (10, 20, 30, 40, 50) is compared. According to Tables 26 and 27, our method consistently outperforms the baseline across all eight testing tasks and performs well even with as few as 10 shots. Due to limitations in the number of images generated by DALL·E-3, we train the model only once for each testing task.

L Model Attribution with Multiple Source Models

Since this involves multi-source origin attribution, the evaluation metric used in this section is **accuracy** instead of AUC. We evaluate the model under the following conditions: 2 classes (SD, VQ-D), 4 classes (SD, VQ-D, LDM, Glide), 6 classes (SD, VQ-D, LDM, Glide, GALIP, ProGAN), and 8 classes (SD, VQ-D, LDM, Glide, GALIP, ProGAN, StyleGAN2, GauGAN). Table 28 shows the accuracy of verifying multiple sources using the eight trained one-class classifiers. Table 29 shows the accuracy of directly training a multi-class classifier.

M Comparison with Related Work

There are two other related works. However, they have different settings and use methods distinct from ours.

[53] proposed a method called RONAN (Reverse-engineering-based Origin Attribution) for origin attribution. The authors aimed to create an alteration-free method that does not impair generation quality and is suitable for pre-trained models. RONAN is based on a reverse-engineering algorithm applied to a source model, utilizing the reconstruction loss of reverse-engineering for inference. This process requires access to the source model and generates many additional images for auxiliary determination. For instance, to determine if an image is generated by a specific model, 100 images need to be generated for inference via Grubbs' Hypothesis Testing [12].

[22] proposed a method called FLIPAD (Final-Layer Inversion Plus Anomaly Detection) to determine if a sample is generated by a specific model. The authors aimed to develop a method not limited to a closed-world setting and without the need for undesirable changes to the target model. FLIPAD is based on final-layer inversion for feature extraction and anomaly detection for inference. This process requires access to the source model and needs 10,000 samples during training.

In contrast, we aim to conduct origin attribution in an open world with only access to a few images generated by a source model. We formulate this problem as a few-shot one-class classification task. Our method, OCC-CLIP, is based on the pre-trained CLIP model and utilizes prompt engineering and adversarial data augmentation. It does not require access to any model and does not need to generate extra images for auxiliary determination. In the training phase, we only need 50, or even 10, samples from both non-target and target image datasets to perform effectively.

N Prompts for DALL·E-3

The following prompts were randomly selected from the annotations of the validation set of the COCO dataset.

A man holding a motion controlled video game controller
A person sitting at a table filled with mexican food.
a work desk with a monitor and keyboard
A clock is standing in the middle of the grass in the middle of the afternoon.
some people are walking in front of a tall building
Three people standing before airport counters below airport signs.
A snow mountain being used for winter sports.
a man that is cutting up some kind of fruit
A bench sitting on a sidewalk near a line of cars.
A skier holds his skis as he stands in the mountains.
A living room filled with lots of furniture and seats.

A train is parked at the station loading passengers.
A man and woman with two Clydesdale horses.
A bathroom with tub and toilet, tiled in white tiling.
A Amtrak train traveling on a railroad tracks.
A stuffed teddy bear sitting amongst pillows on a bed
A green train traveling down train tracks next to another train.
A set of three red double decker buses parked next to each other.
A large plate of waffles is next to plates with slices of peaches.
A couple of people riding skis down a snow covered slope.
A big building perched atop a hill with a sign in the foreground.
The man and woman are playing video games in the room.
Barbecued meat and vegetables laid out on a counter ready for dinner
A herd of elephants in the wild near a river.
A hot dog with mustard and a bun next to a ketchup cup.
A dog rides on the back of a sheep.
The woman is holding up two large hot dogs.
A person is flying through the air near some mountains on his snowboard.
A jockey riding and jumping with a horse in an obstacle course.
a big bed with a lamp and bedside table and sliding glass leading to a balcony
Several lambs and sheep standing on hay and eating it.
an image of a statue of men in a carriage being driven by horses
There is a hose hooked up to the fire hydrant by the building.
A FOUNTAIN IN A ROUNDABOUT WITH PEOPLE PASSING BY
A heard of sheep together in a wheat field.
A woman opening a suitcase on the bed
A traffic light over a city street with cars.
A photo taken from behind a fence of tennis players on the court.
A row of parking meters sitting in a park.
A woman is talking on her phone while dragging on a cigarette.
Subway braking on rails in front of metropolitan city
A young man and woman share some pastries.
An aerial photo of a very long train station at night.
a foggy day that has some lights by a road
A dog sitting on the floor between a person legs
This dirt bike rider is smiling and raising his fist in triumph.
A desk top computer and a laptop sitting on a computer desk
an orange brick building with a window and a big mascot
A person doing a trick on a skateboard in the road
LOTS OF CUPCAKES ARRANGED ON A TABLE WITH NAPKINS
A bathroom sink underneath a medicine cabinet next to a window.
A red bus driving in front of a double decker bus.
The girl is about to kick a soccer ball.
A skateboarder jumping through the air and doing a trick.
Woman with surfboard getting kisses from dog at waters edge.
A skier leaning to the side on a snowy hill.

some motorcycles parked and the one in front is a silver three wheeler
A skier standing at the stop of a mountain slope.
Two planes that are flying in the sky.
A man presenting something to another man in a tent.
A man in a wet suit walking with another man in a wet suit with equipment.
a man flying a kite in a big green field
A man holding a dog mug pints the remote.
two little birds sitting on the shore by a piece of wood
A giraffe eats next to a zebra among some rocks.
A cat thats about to take a bute out of someone's sandwich.
A gray dog wearing an orange bow tie laying on a sofa.
A yellow swan boat ride on a body of water under a cloudless sky.
A old house with a clock tower in brown and white.
a surfer that has fallen off of his surf board
Someone is enjoying a small slice of pie.
A coach balances a soccer ball in front of her team
A couple of giraffes eating out of a feeding bin in a zoo type facility.
A large gray teddy bear sitting next to a candle.
an image of a cat lying next to a stuffed animal
A woman smiling while she prepares a plate of food.
A group of zebras are standing in a desert.
Cat sitting on a window sill in front of a windmill.
A bunch of people walking on the street with umbrellas.
A skate boarder doing jumps at night on city street.
A man with glasses holding a baby in a white outfit.
A man in glasses talking on a cellphone.
A man with glasses is smiling and clapping.
two people sitting on benches with trees in the background
A girl riding on the back of a scooter on a cobbled road.
A small computer keyboard, box and matching mouse
a couple of people on skis stand in the snow
A man sitting in a chair with an infant laying on his lap and a woman standing
over the top and looking down at the infant.
Several alcoholic beverages and mixers on a airplane tray.
A very big pretty horse pulling a fancy carriage.
A young man has just hit a baseball with the bat
A woman in a green shirt stands near a table with bowls or oranges on it in a
market.
A person watching a small plane taking off in a field
A bicycle being held by a bar on a train.
A baseball game in progress in the open air.
A police car parked on the side of the road.
A white toilet sitting next to a toilet paper roller.
A bright green frog on a bright green plant.
A man surfing a wave on his surf board.

A small dog reclines on a piece of furniture next to an electronic game controller.

A person skiing alone in the snow capped area

A large family group at a round table in a restaraunt.

A stop sign is erected next to a flower bush.

Table 10: Evaluation of OCC-CLIP, along with 14 basic methods. Training samples are sourced from SD (denoted as target) and COCO (denoted as non-target). The test data are assembled from SD as well as a different generative image dataset. The optimal outcomes for individual datasets are emphasized using **bold** formatting.

Methods	VQ-D	LDM	Glide	GALIP
VGG16 [44]	0.6458 \pm 2.64e-2	0.5438 \pm 2.29e-2	0.5652 \pm 2.48e-2	0.7017 \pm 4.15e-2
ResNet50 [14]	0.6693 \pm 5.27e-2	0.5485 \pm 3.07e-2	0.6103 \pm 4.03e-2	0.7307 \pm 3.94e-2
Inception-v3 [47]	0.6378 \pm 4.45e-2	0.5349 \pm 2.49e-2	0.5252 \pm 1.78e-2	0.6759 \pm 4.26e-2
DenseNet-121 [16]	0.7264 \pm 4.70e-2	0.5337 \pm 3.40e-2	0.5730 \pm 5.62e-2	0.7478 \pm 4.85e-2
ViT-B-16 [9]	0.7502 \pm 2.70e-2	0.5387 \pm 2.06e-2	0.6069 \pm 4.46e-2	0.7837 \pm 2.05e-2
DeiT-B-16 [50]	0.6825 \pm 3.96e-2	0.5366 \pm 2.47e-2	0.5698 \pm 2.91e-2	0.6592 \pm 2.76e-2
CaiT-S-24 [51]	0.6522 \pm 4.85e-2	0.5301 \pm 2.51e-2	0.5550 \pm 5.18e-2	0.6466 \pm 4.48e-2
Swin-B-4 [25]	0.8634 \pm 3.69e-2	0.7180 \pm 5.06e-2	0.7473 \pm 3.97e-2	0.8478 \pm 4.14e-2
Image_Patch [28]	0.6638 \pm 3.69e-2	0.5269 \pm 1.84e-2	0.5534 \pm 3.15e-2	0.6847 \pm 4.41e-2
Feature_Patch [28]	0.6999 \pm 2.53e-2	0.6116 \pm 3.27e-2	0.7385 \pm 7.70e-2	0.6285 \pm 5.62e-2
CLIP [33]	0.7272 \pm 0.0	0.7243 \pm 0.0	0.6078 \pm 0.0	0.7304 \pm 0.0
CoOp [59]	0.9503 \pm 1.68e-2	0.8373 \pm 5.33e-2	0.9266 \pm 3.19e-2	0.9660 \pm 2.56e-2
OCC-CLIP	0.9703\pm9.06e-3	0.8801\pm2.06e-2	0.9519\pm1.45e-2	0.9798\pm1.18e-2

Table 11: Evaluation of OCC-CLIP, along with 14 basic methods. Training samples are sourced from SD (denoted as target) and COCO (denoted as non-target). The test data are assembled from SD as well as a different generative image dataset. The optimal outcomes for individual datasets are emphasized using **bold** formatting.

Methods	ProGan	StyleGan2	GauGan	Overall
VGG16 [44]	0.6609 \pm 5.46e-2	0.6170 \pm 7.87e-2	0.7096 \pm 4.62e-2	0.6349 \pm 4.61e-2
ResNet50 [14]	0.6646 \pm 4.56e-2	0.6458 \pm 9.33e-2	0.7475 \pm 3.39e-2	0.6595 \pm 5.18e-2
Inception-v3 [47]	0.6453 \pm 4.87e-2	0.5947 \pm 4.88e-2	0.6982 \pm 4.31e-2	0.6160 \pm 4.03e-2
DenseNet-121 [16]	0.7204 \pm 3.35e-2	0.6645 \pm 6.78e-2	0.8091 \pm 2.94e-2	0.6821 \pm 4.70e-2
ViT-B-16 [9]	0.6062 \pm 3.71e-2	0.6593 \pm 6.11e-2	0.7259 \pm 2.77e-2	0.6673 \pm 3.67e-2
DeiT-B-16 [50]	0.5721 \pm 2.77e-2	0.5898 \pm 4.20e-2	0.6507 \pm 2.72e-2	0.6087 \pm 3.17e-2
CaiT-S-24 [51]	0.5697 \pm 2.63e-2	0.5804 \pm 5.12e-2	0.6645 \pm 3.42e-2	0.5998 \pm 4.17e-2
Swin-B-4 [25]	0.5879 \pm 3.68e-2	0.7054 \pm 6.30e-2	0.7822 \pm 3.91e-2	0.7503 \pm 4.48e-2
Image_Patch [28]	0.6624 \pm 4.82e-2	0.6706 \pm 9.62e-2	0.7674 \pm 4.09e-2	0.6470 \pm 5.05e-2
Feature_Patch [28]	0.7136 \pm 2.77e-2	0.8546 \pm 2.65e-2	0.8152 \pm 3.14e-2	0.7231 \pm 4.35e-2
CLIP [33]	0.5229 \pm 0.0	0.5393 \pm 0.0	0.6463 \pm 0.0	0.6426 \pm 0.0
CoOp [59]	0.8861 \pm 4.46e-2	0.9533 \pm 2.30e-2	0.9643 \pm 2.91e-2	0.9263 \pm 3.41e-2
OCC-CLIP	0.9452\pm3.14e-2	0.9651\pm1.54e-2	0.9910\pm7.32e-3	0.9548\pm1.75e-2

Table 12: Evaluation sensitivity of OCC-CLIP to Source Models and choice of Epsilon. The leftmost column represents target datasets. In training phase, the non-target images are from COCO. In testing phase, the non-target images are from a different generative image dataset shown in the first row. Different epsilon is chosen: 0.0, 0.03125, 0.0625, 0.1, 0.2, and 0.5. The optimal outcomes for individual datasets are emphasized using **bold** formatting.

Target	Epsilon	SD _{LAION-5B}	VQ-D _{coco}	LDM _{LAION-400M}	Glide _{filtered-cc}
SD _{LAION-5B}	0.0	–	0.9503±1.68e-2	0.8373±5.33e-2	0.9266±3.19e-2
	0.03125	–	0.9627±1.74e-2	0.8782±4.48e-2	0.9248±3.58e-2
	0.0625	–	0.9644±1.73e-2	0.8720±3.74e-2	0.9332±3.10e-2
	0.1	–	0.9703±9.06e-3	0.8801±2.06e-2	0.9519±1.45e-2
	0.2	–	0.9595±1.49e-2	0.8456±4.05e-2	0.9388±2.54e-2
	0.5	–	0.9601±1.76e-2	0.8475±4.87e-2	0.9384±3.18e-2
VQ-D _{coco}	0.0	0.9988±9.79e-4	–	0.7337±4.63e-2	0.7435±4.39e-2
	0.03125	0.9982±5.66e-4	–	0.7077±4.06e-2	0.6988±4.61e-2
	0.0625	0.9986±1.08e-3	–	0.6783±5.50e-2	0.7112±5.12e-2
	0.1	0.9992±6.54e-4	–	0.6924±6.56e-2	0.7264±5.62e-2
	0.2	0.9992±3.85e-4	–	0.6998±3.20e-2	0.7478±4.84e-2
	0.5	0.9994±3.74e-4	–	0.6887±5.88e-2	0.7468±4.32e-2
LDM _{LAION-400M}	0.0	0.9925±7.39e-3	0.6793±6.63e-2	–	0.6468±6.15e-2
	0.03125	0.9941±4.49e-3	0.7295±3.77e-2	–	0.6502±6.56e-2
	0.0625	0.9945±5.66e-3	0.7473±4.92e-2	–	0.6761±5.30e-2
	0.1	0.9957±3.97e-3	0.7507±5.51e-2	–	0.6847±5.32e-2
	0.2	0.9958±2.73e-3	0.7461±5.98e-2	–	0.6980±6.25e-2
	0.5	0.9951±3.13e-3	0.7404±5.12e-2	–	0.7142±3.77e-2
Glide _{filtered-cc}	0.0	0.9985±1.86e-3	0.8573±4.27e-2	0.8314±3.08e-2	–
	0.03125	0.9997±3.16e-4	0.8970±1.53e-2	0.8553±2.96e-2	–
	0.0625	0.9998±1.90e-4	0.8919±1.76e-2	0.8559±2.49e-2	–
	0.1	0.9998±2.47e-4	0.8958±2.18e-2	0.8585±2.50e-2	–
	0.2	0.9999±1.18e-4	0.9073±1.04e-2	0.8731±2.42e-2	–
	0.5	0.9999±1.08e-4	0.8898±2.37e-2	0.8724±3.14e-2	–
GALIP _{coco}	0.0	0.9999±2.29e-4	0.9036±2.75e-2	0.8441±4.62e-2	0.7802±4.74e-2
	0.03125	0.9996±1.85e-4	0.8970±4.26e-2	0.8256±3.54e-2	0.7044±3.36e-2
	0.0625	0.9998±1.83e-4	0.9205±2.30e-2	0.8165±3.54e-2	0.7545±4.19e-2
	0.1	0.9999±8.94e-5	0.9345±1.99e-2	0.8626±1.99e-2	0.7779±2.88e-2
	0.2	0.9999±6.40e-5	0.9380±1.81e-2	0.8776±2.96e-2	0.8238±2.09e-2
	0.5	1.0000±4.00e-5	0.9332±1.57e-2	0.8907±3.24e-2	0.8361±1.73e-2
ProGan _{lsun}	0.0	0.9972±1.91e-3	0.9475±2.01e-2	0.9544±1.86e-2	0.9453±2.00e-2
	0.03125	0.9904±6.10e-3	0.9429±2.49e-2	0.9646±2.12e-2	0.8852±5.46e-2
	0.0625	0.9934±4.67e-3	0.9313±3.26e-2	0.9517±2.18e-2	0.8868±4.72e-2
	0.1	0.9961±3.27e-3	0.9477±2.46e-2	0.9585±2.58e-2	0.9320±3.59e-2
	0.2	0.9970±1.73e-3	0.9689±7.48e-3	0.9762±7.41e-3	0.9639±1.40e-2
	0.5	0.9980±1.50e-3	0.9749±1.13e-2	0.9779±8.49e-3	0.9750±1.16e-2
StyleGan2 _{lsun}	0.0	0.9992±8.01e-4	0.9856±8.77e-3	0.9850±1.12e-2	0.9584±3.33e-2
	0.03125	0.9989±8.72e-4	0.9921±6.86e-3	0.9894±1.14e-2	0.9352±3.12e-2
	0.0625	0.9995±3.75e-4	0.9911±9.81e-3	0.9877±1.32e-2	0.9511±2.88e-2
	0.1	0.9996±4.63e-4	0.9937±5.33e-3	0.9851±2.31e-2	0.9652±2.42e-2
	0.2	0.9997±2.32e-4	0.9941±2.36e-3	0.9900±4.43e-3	0.9677±1.23e-2
	0.5	0.9997±1.66e-4	0.9943±2.17e-3	0.9917±3.25e-3	0.9688±9.89e-3
GauGan _{coco}	0.0	0.9991±6.97e-4	0.9388±2.52e-2	0.9862±6.91e-3	0.9757±1.02e-2
	0.03125	0.9960±2.16e-3	0.8420±6.06e-2	0.9616±1.45e-2	0.8890±4.38e-2
	0.0625	0.9976±1.32e-3	0.8593±3.97e-2	0.9678±8.16e-3	0.9149±2.89e-2
	0.1	0.9982±1.33e-3	0.9132±2.93e-2	0.9745±9.90e-3	0.9593±2.54e-2
	0.2	0.9992±1.34e-3	0.9456±2.44e-2	0.9881±5.31e-3	0.9904±3.47e-3
	0.5	0.9997±2.94e-4	0.9518±2.79e-2	0.9908±5.50e-3	0.9914±4.66e-3

Table 13: Evaluation sensitivity of OCC-CLIP to Source Models and choice of Epsilon. The leftmost column represents target datasets. In training phase, the non-target images are from COCO. In testing phase, the non-target images are from a different generative image dataset shown in the first row. Different epsilon is chosen: 0.0, 0.03125, 0.0625, 0.1, 0.2, and 0.5. The optimal outcomes for individual datasets are emphasized using **bold** formatting.

Target	Epsilon	GALIP _{coco}	ProGan _{lsun}	StyleGan2 _{lsun}	GauGan _{coco}	Overall
SD _{LAION-5B}	0.0	0.9660±2.56e-2	0.8861±4.46e-2	0.9533±2.30e-2	0.9643±2.91e-2	0.9263±3.41e-2
	0.03125	0.9764±1.83e-2	0.9564±1.84e-2	0.9695±2.23e-2	0.9668±2.79e-3	0.9521±2.61e-2
	0.0625	0.9762±1.77e-2	0.9470±3.82e-2	0.9684±1.67e-2	0.9939±7.33e-3	0.9507±2.61e-2
	0.1	0.9798±1.18e-2	0.9452±3.14e-2	0.9651±1.54e-2	0.9910±7.32e-3	0.9548±1.75e-2
	0.2	0.9722±1.38e-2	0.9196±2.95e-2	0.9578±2.46e-2	0.9821±1.04e-2	0.9394±2.47e-2
	0.5	0.9670±2.01e-2	0.9002±3.15e-2	0.9627±1.99e-2	0.9816±1.03e-2	0.9368±2.83e-2
VQ-D _{coco}	0.0	0.7558±3.82e-2	0.9290±2.80e-2	0.9924±3.87e-3	0.9352±1.95e-2	0.8698±3.09e-2
	0.03125	0.7303±5.59e-2	0.9920±4.31e-3	0.9971±2.20e-3	0.9946±2.51e-3	0.8741±3.15e-2
	0.0625	0.7245±6.14e-2	0.9923±5.29e-3	0.9938±5.97e-3	0.9949±3.19e-3	0.8705±3.68e-2
	0.1	0.7327±5.82e-2	0.9931±3.58e-3	0.9936±5.51e-3	0.9936±2.50e-3	0.8758±3.95e-2
	0.2	0.7120±4.83e-2	0.9938±2.40e-3	0.9958±1.95e-3	0.9931±2.39e-3	0.8774±2.86e-2
	0.5	0.6876±5.09e-2	0.9899±4.30e-3	0.9960±2.28e-3	0.9909±3.81e-3	0.8713±3.37e-2
LDM _{LAION-400M}	0.0	0.6263±7.51e-2	0.9565±4.27e-2	0.9896±8.87e-3	0.9758±2.86e-2	0.8381±4.87e-2
	0.03125	0.6276±6.30e-2	0.9970±2.09e-3	0.9957±3.08e-3	0.9997±3.55e-4	0.8563±3.73e-2
	0.0625	0.6389±4.84e-2	0.9971±1.94e-3	0.9957±3.73e-3	0.9997±4.27e-4	0.8642±3.30e-2
	0.1	0.6530±4.87e-2	0.9956±3.43e-3	0.9940±3.47e-3	0.9992±7.47e-4	0.8676±3.44e-2
	0.2	0.6515±5.97e-2	0.9919±3.47e-3	0.9939±2.97e-3	0.9977±1.57e-3	0.8679±3.98e-2
	0.5	0.6388±3.67e-2	0.9905±4.88e-3	0.9952±2.17e-3	0.9973±1.68e-3	0.8674±2.79e-2
Glide _{filtered-cc}	0.0	0.6687±8.82e-2	0.9629±4.83e-2	0.9916±6.77e-3	0.9814±3.92e-2	0.8988±4.55e-2
	0.03125	0.6910±4.59e-2	0.9974±2.15e-3	0.9980±1.48e-3	0.9998±4.35e-4	0.9197±2.15e-2
	0.0625	0.6942±7.30e-2	0.9973±2.94e-3	0.9951±3.15e-3	0.9996±7.15e-4	0.9191±2.99e-2
	0.1	0.6834±6.43e-2	0.9974±1.66e-3	0.9949±2.73e-3	0.9997±4.78e-4	0.9185±2.74e-2
	0.2	0.6406±4.52e-2	0.9965±1.45e-3	0.9976±1.20e-3	0.9998±2.07e-4	0.9164±1.98e-2
	0.5	0.6535±6.54e-2	0.9926±6.24e-3	0.9961±2.96e-3	0.9989±1.91e-3	0.9147±2.90e-2
GALIP _{coco}	0.0	-	0.9982±1.16e-3	0.9992±7.43e-4	0.9996±5.71e-4	0.9321±2.71e-2
	0.03125	-	0.9999±1.10e-4	0.9994±5.35e-4	1.0000±0.0	0.9180±2.45e-2
	0.0625	-	0.9999±7.00e-5	0.9994±5.14e-4	1.0000±0.0	0.9272±2.25e-2
	0.1	-	0.9999±1.35e-4	0.9993±5.58e-4	1.0000±0.0	0.9392±1.52e-2
	0.2	-	0.9999±8.06e-5	0.9995±4.41e-4	1.0000±0.0	0.9484±1.53e-2
	0.5	-	0.9997±2.84e-4	0.9997±2.68e-4	1.0000±3.00e-5	0.9513±1.51e-2
ProGan _{lsun}	0.0	0.9818±1.24e-2	-	0.9278±1.49e-2	0.7993±2.00e-2	0.9362±1.66e-2
	0.03125	0.9856±9.97e-3	-	0.9065±3.16e-2	0.9249±1.62e-2	0.9429±2.79e-2
	0.0625	0.9830±1.18e-2	-	0.8478±3.35e-2	0.9146±1.04e-2	0.9298±2.71e-2
	0.1	0.9885±7.41e-3	-	0.8471±4.24e-2	0.8885±1.16e-2	0.9369±2.55e-2
	0.2	0.9937±3.71e-3	-	0.8719±3.09e-2	0.8456±2.75e-2	0.9453±1.71e-2
	0.5	0.9936±4.05e-3	-	0.9109±3.60e-2	0.8698±2.34e-2	0.9572±1.77e-2
StyleGan2 _{lsun}	0.0	0.9849±1.35e-2	0.8687±3.78e-2	-	0.9543±1.79e-2	0.9623±2.15e-2
	0.03125	0.9925±4.85e-3	0.9585±2.45e-2	-	0.9952±2.71e-3	0.9803±1.60e-2
	0.0625	0.9917±5.44e-3	0.9685±1.35e-2	-	0.9953±2.49e-3	0.9836±1.37e-2
	0.1	0.9926±6.85e-3	0.9649±1.50e-2	-	0.9946±1.80e-3	0.9851±1.43e-2
	0.2	0.9924±3.55e-3	0.9500±1.69e-2	-	0.9889±4.87e-3	0.9833±8.43e-3
	0.5	0.9912±4.58e-3	0.9407±2.38e-2	-	0.9885±4.97e-3	0.9821±1.02e-2
GauGan _{coco}	0.0	0.9901±6.44e-3	0.6812±6.34e-2	0.9676±1.29e-2	-	0.9368±2.68e-2
	0.03125	0.9862±1.03e-2	0.6965±5.18e-2	0.9225±3.06e-2	-	0.8991±3.69e-2
	0.0625	0.9908±4.86e-3	0.7132±6.78e-2	0.8919±3.83e-2	-	0.9051±3.50e-2
	0.1	0.9958±3.90e-3	0.7752±6.78e-2	0.9425±3.78e-2	-	0.9370±3.31e-2
	0.2	0.9980±1.97e-3	0.8039±4.45e-2	0.9492±1.26e-2	-	0.9535±1.99e-2
	0.5	0.9976±2.05e-3	0.7709±4.78e-2	0.9583±1.30e-2	-	0.9515±2.17e-2

Table 14: Evaluation sensitivity of different ways of doing ADA on different open world real image datasets. 1) In training phase, conditioned on doing ADA only on half of non-target image set, one of COCO, ImageNet, Flickr, or CC12M is used as the non-target dataset. 2) Conditioned on using non-target images selected from COCO, we do ADA on half of target image set, on half of the both non-target images and target images, and on half of the target images which are then treated as non-target ones.

Methods	VQ-D	LDM	Glide	GALIP
COCO+Neg	0.9703 \pm 9.06e-3	0.8801 \pm 2.06e-2	0.9519 \pm 1.45e-2	0.9798\pm1.18e-2
Flickr+Neg	0.9586 \pm 2.17e-2	0.8775 \pm 2.86e-2	0.8805 \pm 4.33e-2	0.9546 \pm 1.89e-2
CC12M+Neg	0.8995 \pm 2.68e-2	0.8523 \pm 3.96e-2	0.9364 \pm 1.74e-2	0.9444 \pm 2.73e-2
ImageNet+Neg	0.9773\pm7.67e-3	0.9136\pm2.51e-2	0.9664\pm1.39e-2	0.9710 \pm 2.34e-2
COCO+None	0.9503 \pm 1.68e-2	0.8373 \pm 5.33e-2	0.9266 \pm 3.19e-2	0.9660 \pm 2.56e-2
COCO+Neg	0.9703\pm9.06e-3	0.8801 \pm 2.06e-2	0.9519 \pm 1.45e-2	0.9798 \pm 1.18e-2
COCO+Both	0.9346 \pm 2.32e-2	0.8229 \pm 5.47e-2	0.9070 \pm 2.81e-2	0.9655 \pm 1.86e-2
COCO+Target	0.8779 \pm 6.93e-2	0.7618 \pm 8.77e-2	0.8728 \pm 4.90e-2	0.9205 \pm 3.83e-2
COCO+T-NT	0.9684 \pm 2.45e-2	0.8818\pm4.91e-2	0.9522\pm2.98e-2	0.9803\pm1.26e-2

Table 15: Evaluation sensitivity of different ways of doing ADA on different open world real image datasets. 1) In training phase, conditioned on doing ADA only on half of non-target image set, one of COCO, ImageNet, Flickr, or CC12M is used as the non-target dataset. 2) Conditioned on using non-target images selected from COCO, we do ADA on half of target image set, on half of the both non-target images and target images, and on half of the target images which are then treated as non-target ones.

Methods	ProGan	StyleGan2	GauGan	Overall
COCO+Neg	0.9452 \pm 3.14e-2	0.9651 \pm 1.54e-2	0.9910 \pm 7.32e-3	0.9548 \pm 1.75e-2
Flickr+Neg	0.9494 \pm 2.17e-2	0.9368 \pm 3.01e-2	0.9908 \pm 9.65e-3	0.9355 \pm 2.67e-2
CC12M+Neg	0.9671 \pm 1.32e-2	0.9422 \pm 2.45e-2	0.9822 \pm 6.75e-3	0.9320 \pm 2.43e-2
ImageNet+Neg	0.9862\pm4.11e-3	0.9861\pm4.22e-3	0.9966\pm1.16e-3	0.9710\pm1.44e-2
COCO+None	0.8861 \pm 4.46e-2	0.9533 \pm 2.30e-2	0.9643 \pm 2.91e-2	0.9263 \pm 3.41e-2
COCO+Neg	0.9452\pm3.14e-2	0.9651 \pm 1.54e-2	0.9910\pm7.32e-3	0.9548\pm1.75e-2
COCO+Both	0.8637 \pm 5.45e-2	0.9113 \pm 5.48e-2	0.9511 \pm 2.26e-2	0.9080 \pm 3.99e-2
COCO+Target	0.6470 \pm 1.15e-1	0.9302 \pm 3.77e-2	0.7553 \pm 1.13e-1	0.8236 \pm 7.90e-2
COCO+T-NT	0.9244 \pm 4.53e-2	0.9665\pm2.08e-2	0.9859 \pm 1.18e-2	0.9514 \pm 3.09e-2

Table 16: Evaluation of OCC-CLIP Relative to the Proportion of Augmented Non-Target Images. Five proportions are investigated: 0%, 25%, 50%, 75%, and 100%.

Methods	VQ-D	LDM	Glide	GALIP
OCC-CLIP 0%	0.9503 \pm 1.68e-2	0.8373 \pm 5.33e-2	0.9266 \pm 3.19e-2	0.9660 \pm 2.56e-2
OCC-CLIP 25%	0.9558 \pm 2.19e-2	0.8519 \pm 4.65e-2	0.9356 \pm 3.23e-2	0.9754 \pm 1.55e-2
OCC-CLIP 50%	0.9703\pm9.06e-3	0.8801 \pm 2.06e-2	0.9519 \pm 1.45e-2	0.9798 \pm 1.18e-2
OCC-CLIP 75%	0.9683 \pm 1.25e-2	0.8828\pm3.34e-2	0.9531\pm2.17e-2	0.9815\pm9.53e-3
OCC-CLIP 100%	0.9252 \pm 1.72e-2	0.8414 \pm 2.19e-2	0.8996 \pm 2.84e-2	0.9614 \pm 8.53e-3

Table 17: Evaluation of OCC-CLIP Relative to the Proportion of Augmented Non-Target Images. Five proportions are investigated: 0%, 25%, 50%, 75%, and 100%.

Methods	ProGan	StyleGan2	GauGan	Overall
OCC-CLIP 0%	0.8861 \pm 4.46e-2	0.9533 \pm 2.30e-2	0.9643 \pm 2.91e-2	0.9263 \pm 3.41e-2
OCC-CLIP 25%	0.9209 \pm 5.46e-2	0.9561 \pm 3.20e-2	0.9878 \pm 1.15e-2	0.9405 \pm 3.39e-2
OCC-CLIP 50%	0.9452 \pm 3.14e-2	0.9651 \pm 1.54e-2	0.9910 \pm 7.32e-3	0.9548 \pm 1.75e-2
OCC-CLIP 75%	0.9466\pm2.23e-2	0.9675\pm1.58e-2	0.9911\pm5.11e-3	0.9558\pm1.93e-2
OCC-CLIP 100%	0.9393 \pm 1.86e-2	0.8674 \pm 2.64e-2	0.9861 \pm 5.64e-3	0.9172 \pm 1.98e-2

Table 18: Evaluation of OCC-CLIP on Different Numbers of Shots. This table shows the mean origin attribution performance of CoOp and OCC-CLIP on 7 different testing tasks with a variable number of shots: 10, 20, 30, 40, 50, 100, and 200. The target dataset is SD, while the non-target dataset is from one of COCO, CC12M, Flickr, or ImageNet.

Non-Target	Methods	10	20	30	40
COCO	CoOp	0.7771 \pm 1.02e-1	0.8500 \pm 8.36e-2	0.8836 \pm 4.45e-2	0.9183 \pm 4.57e-2
	OCC-CLIP	0.8842\pm7.22e-2	0.9293\pm4.45e-2	0.9435\pm2.23e-2	0.9494\pm2.62e-2
Flickr	CoOp	0.7905 \pm 7.82e-2	0.8649 \pm 7.54e-2	0.8784 \pm 4.39e-2	0.8734 \pm 8.64e-2
	OCC-CLIP	0.8921\pm8.26e-2	0.9231\pm3.86e-2	0.9360\pm2.32e-2	0.9358\pm3.88e-2
CC12M	CoOp	0.7737 \pm 8.18e-2	0.8302 \pm 7.79e-2	0.8640 \pm 5.67e-2	0.8582 \pm 4.77e-2
	OCC-CLIP	0.8651\pm8.11e-2	0.9103\pm4.89e-2	0.9169\pm3.68e-2	0.9234\pm4.18e-2
ImageNet	CoOp	0.8559 \pm 8.04e-2	0.8938 \pm 5.05e-2	0.9308 \pm 3.49e-2	0.9299 \pm 4.02e-2
	OCC-CLIP	0.9192\pm6.34e-2	0.9489\pm2.87e-2	0.9647\pm1.86e-2	0.9673\pm2.15e-2

Table 19: Evaluation of OCC-CLIP on Different Numbers of Shots. This table shows the mean origin attribution performance of CoOp and OCC-CLIP on 7 different testing tasks with a variable number of shots: 10, 20, 30, 40, 50, 100, and 200. The target dataset is SD, while the non-target dataset is from one of COCO, CC12M, Flickr, or ImageNet.

Non-Target	Methods	50	100	200
COCO	CoOp	0.9263 \pm 3.41e-2	0.9541 \pm 2.21e-2	0.9601 \pm 2.00e-2
	OCC-CLIP	0.9548 \pm 1.75e-2	0.9684 \pm 1.56e-2	0.9746 \pm 1.08e-2
Flickr	CoOp	0.8788 \pm 7.17e-2	0.9071 \pm 3.68e-2	0.9177 \pm 3.46e-2
	OCC-CLIP	0.9355 \pm 2.67e-2	0.9531 \pm 2.88e-2	0.9588 \pm 3.35e-2
CC12M	CoOp	0.8689 \pm 4.27e-2	0.8973 \pm 3.62e-2	0.8922 \pm 3.70e-2
	OCC-CLIP	0.9320 \pm 2.43e-2	0.9342 \pm 3.18e-2	0.9331 \pm 3.55e-2
ImageNet	CoOp	0.9377 \pm 3.12e-2	0.9507 \pm 2.28e-2	0.9663 \pm 1.57e-2
	OCC-CLIP	0.9710 \pm 1.44e-2	0.9768 \pm 5.42e-3	0.9813 \pm 8.37e-3

Table 20: Evaluation of OCC-CLIP under Image Processing. Six image processing methods are executed: Gaussian Blur, Gaussian Noise, Grayscale, Rotation, Flip, and a mixture of all the these data augmentation methods.

Process	Method	VQ-D	LDM	Glide	GALIP
None	CoOp	0.9503 \pm 1.68e-2	0.8373 \pm 5.33e-2	0.9266 \pm 3.19e-2	0.9660 \pm 2.56e-2
	OCC-CLIP	0.9703 \pm 9.06e-3	0.8801 \pm 2.06e-2	0.9519 \pm 1.45e-2	0.9798 \pm 1.18e-2
Gaussian Blur	CoOp	0.8918 \pm 3.15e-2	0.7067 \pm 6.27e-2	0.8555 \pm 3.86e-2	0.9240 \pm 4.39e-2
	OCC-CLIP	0.9286 \pm 2.54e-2	0.7700 \pm 4.71e-2	0.8914 \pm 3.59e-2	0.9502 \pm 2.53e-2
Gaussian Noise	CoOp	0.9159 \pm 3.99e-2	0.7646 \pm 8.89e-2	0.8847 \pm 5.81e-2	0.9374 \pm 4.72e-2
	OCC-CLIP	0.9348 \pm 2.60e-2	0.7877 \pm 6.07e-2	0.8993 \pm 4.49e-2	0.9479 \pm 4.08e-2
Grayscale	CoOp	0.8455 \pm 5.94e-2	0.6348 \pm 8.15e-2	0.7987 \pm 7.14e-2	0.8852 \pm 8.56e-2
	OCC-CLIP	0.8817 \pm 4.82e-2	0.6575 \pm 9.85e-2	0.8254 \pm 5.91e-2	0.9085 \pm 5.80e-2
Rotate 90	CoOp	0.9319 \pm 2.56e-2	0.7996 \pm 6.96e-2	0.9033 \pm 4.44e-2	0.9514 \pm 3.88e-2
	OCC-CLIP	0.9547 \pm 1.29e-2	0.8406 \pm 2.94e-2	0.9290 \pm 2.07e-2	0.9680 \pm 1.78e-2
Flip	CoOp	0.9481 \pm 2.23e-2	0.8327 \pm 6.56e-2	0.9233 \pm 3.89e-2	0.9634 \pm 3.23e-2
	OCC-CLIP	0.9650 \pm 1.27e-2	0.8675 \pm 3.28e-2	0.9434 \pm 2.39e-2	0.9754 \pm 1.56e-2
Mixture	CoOp	0.7858 \pm 6.29e-2	0.5855 \pm 6.07e-2	0.7314 \pm 8.42e-2	0.8405 \pm 8.70e-2
	OCC-CLIP	0.7895 \pm 6.61e-2	0.5959 \pm 4.79e-2	0.7137 \pm 7.38e-2	0.8276 \pm 8.91e-2

Table 21: Evaluation of OCC-CLIP under Image Processing. Six image processing methods are executed: Gaussian Blur, Gaussian Noise, Grayscale, Rotation, Flip, and a mixture of all the these data augmentation methods.

Process	Method	ProGan	StyleGan2	GauGan	Overall
None	CoOp	0.8861±4.46e-2	0.9533±2.30e-2	0.9643±2.91e-2	0.9263±3.41e-2
	OCC-CLIP	0.9452±3.14e-2	0.9651±1.54e-2	0.9910±7.32e-3	0.9548±1.75e-2
Gaussian Blur	CoOp	0.7772±8.53e-2	0.9042±3.74e-2	0.9181±6.47e-2	0.8539±5.50e-2
	OCC-CLIP	0.8796±4.99e-2	0.9072±4.48e-2	0.9743±1.54e-2	0.9002±3.69e-2
Gaussian Noise	CoOp	0.8298±6.86e-2	0.9269±3.56e-2	0.9418±5.39e-2	0.8859±5.85e-2
	OCC-CLIP	0.8947±3.06e-2	0.9194±3.15e-2	0.9783±1.17e-2	0.9089±3.80e-2
Grayscale	CoOp	0.7002±1.21e-1	0.8636±7.01e-2	0.8839±9.39e-2	0.8017±8.54e-2
	OCC-CLIP	0.8166±1.00e-1	0.8472±8.88e-2	0.9468±6.20e-2	0.8405±7.62e-2
Rotate 90	CoOp	0.8545±5.40e-2	0.9369±3.27e-2	0.9502±3.58e-2	0.9040±4.51e-2
	OCC-CLIP	0.9195±3.33e-2	0.9410±2.49e-2	0.9832±1.12e-2	0.9337±2.28e-2
Flip	CoOp	0.8842±4.61e-2	0.9535±2.55e-2	0.9639±2.98e-2	0.9241±3.97e-2
	OCC-CLIP	0.9384±2.92e-2	0.9567±2.29e-2	0.9892±8.51e-3	0.9479±2.24e-2
Mixture	CoOp	0.6456±9.21e-2	0.7991±8.54e-2	0.8320±1.09e-1	0.7457±8.45e-2
	OCC-CLIP	0.7224±7.87e-2	0.7332±1.16e-1	0.8946±8.37e-2	0.7538±8.17e-2

Table 22: Evaluation sensitivity to choice of prompts. This table shows the performance of OCC-CLIP and CoOp under different choices of label pairs.

Names of Classes	Methods	VQ-D	LDM	Glide	GALIP
fake-real	CoOp	0.9432±2.04e-2	0.8224±5.69e-2	0.9128±3.87e-2	0.9631±1.99e-2
	OCC-CLIP	0.9469±2.19e-2	0.8324±4.15e-2	0.9215±2.33e-2	0.9693±1.58e-2
negative-positive	CoOp	0.9521±2.18e-2	0.8578±4.72e-2	0.9180±4.11e-2	0.9610±2.61e-2
	OCC-CLIP	0.9552±2.46e-2	0.8492±4.55e-2	0.9163±3.34e-2	0.9597±2.22e-2
other-this	CoOp	0.9543±2.99e-2	0.8510±6.57e-2	0.9316±3.63e-2	0.9727±2.51e-2
	OCC-CLIP	0.9611±1.40e-2	0.8587±3.72e-2	0.9202±3.53e-2	0.9712±1.54e-2
real-fake	CoOp	0.9503±1.68e-2	0.8373±5.33e-2	0.9266±3.19e-2	0.9660±2.56e-2
	OCC-CLIP	0.9703±9.06e-3	0.8801±2.06e-2	0.9519±1.45e-2	0.9798±1.18e-2
positive-negative	CoOp	0.9579±1.22e-2	0.8737±3.69e-2	0.9439±1.97e-2	0.9768±1.46e-2
	OCC-CLIP	0.9671±1.50e-2	0.8916±2.60e-2	0.9507±2.56e-2	0.9806±1.35e-2
this-other	CoOp	0.9343±2.91e-2	0.8040±4.93e-2	0.9007±3.66e-2	0.9486±2.31e-2
	OCC-CLIP	0.9488±3.08e-2	0.8438±4.70e-2	0.9268±2.02e-2	0.9654±2.88e-2

Table 23: Evaluation sensitivity to choice of prompts. This table shows the performance of OCC-CLIP and CoOp under different choices of label pairs.

Names of Classes	Methods	ProGan	StyleGan2	GauGan	Overall
fake-real	CoOp	0.8220±5.78e-2	0.9431±2.01e-2	0.9477±2.54e-2	0.9078±3.77e-2
	OCC-CLIP	0.8847±3.59e-2	0.9473±1.90e-2	0.9734±8.03e-3	0.9251±2.60e-2
negative-positive	CoOp	0.8228±7.85e-2	0.9422±3.87e-2	0.9549±2.33e-2	0.9155±4.36e-2
	OCC-CLIP	0.9120±2.99e-2	0.9441±2.86e-2	0.9857±6.60e-3	0.9317±2.94e-2
other-this	CoOp	0.8596±6.79e-2	0.9355±4.04e-2	0.9542±2.79e-2	0.9227±4.50e-2
	OCC-CLIP	0.9096±4.59e-2	0.9501±2.48e-2	0.9791±1.38e-2	0.9357±2.92e-2
real-fake	CoOp	0.8861±4.46e-2	0.9533±2.30e-2	0.9643±2.91e-2	0.9263±3.41e-2
	OCC-CLIP	0.9452±3.14e-2	0.9651±1.54e-2	0.9910±7.32e-3	0.9548±1.75e-2
positive-negative	CoOp	0.8467±6.97e-2	0.9693±1.27e-2	0.9417±3.08e-2	0.9300±3.40e-2
	OCC-CLIP	0.9445±3.51e-2	0.9778±1.98e-2	0.9883±1.18e-2	0.9572±2.24e-2
this-other	CoOp	0.8270±6.63e-2	0.9460±2.65e-2	0.9446±3.26e-2	0.9007±4.02e-2
	OCC-CLIP	0.8932±5.06e-2	0.9559±2.51e-2	0.9766±1.57e-2	0.9301±3.35e-2

Table 24: ADA vs. other data augmentation methods. Different data augmentation methods are used during the training phase.

Methods	VQ-D	LDM	Glide	GALIP
None	0.9503±1.68e-2	0.8373±5.33e-2	0.9266±3.19e-2	0.9660±2.56e-2
Gaussian Blur	0.7780±3.10e-2	0.6417±5.40e-2	0.8460±3.84e-2	0.8034±4.79e-2
Gaussian Noise	0.7673±3.51e-2	0.6078±4.78e-2	0.7735±4.63e-2	0.7004±4.42e-2
Grayscale	0.7490±4.32e-2	0.6111±4.53e-2	0.7794±3.65e-2	0.6877±4.09e-2
Rotate	0.7653±4.36e-2	0.6174±4.94e-2	0.8029±4.12e-2	0.7496±5.03e-2
Flip	0.7614±4.44e-2	0.6182±4.77e-2	0.7880±4.67e-2	0.7370±3.71e-2
Mixture	0.8220±2.72e-2	0.5339±1.49e-2	0.7017±4.14e-2	0.7846±3.55e-2
ADA	0.9703±9.06e-3	0.8801±2.06e-2	0.9519±1.45e-2	0.9798±1.18e-2

Table 25: ADA vs. other data augmentation methods. Different data augmentation methods are used during training phase.

Methods	ProGan	StyleGan2	GauGan	Overall
None	0.8861±4.46e-2	0.9533±2.30e-2	0.9643±2.91e-2	0.9263±3.41e-2
Gaussian Blur	0.9856±1.19e-2	0.9978±1.87e-3	0.9949±5.02e-3	0.8639±3.34e-2
Gaussian Noise	0.9926±5.84e-3	0.9974±2.32e-3	0.9991±1.21e-3	0.8340±3.31e-2
Grayscale	0.9916±4.96e-3	0.9978±1.27e-3	0.9988±9.65e-4	0.8308±3.15e-2
Rotate	0.9930±2.77e-3	0.9985±9.18e-4	0.9985±9.73e-4	0.8465±3.50e-2
Flip	0.9936±3.06e-3	0.9983±1.45e-3	0.9985±1.02e-3	0.8422±3.34e-2
Mixture	0.8743±5.98e-2	0.9328±4.72e-2	0.9440±2.94e-2	0.7991±3.89e-2
ADA	0.9452±3.14e-2	0.9651±1.54e-2	0.9910±7.32e-3	0.9548±1.75e-2

