

Learning Counterfactually Decoupled Attention for Open-World Model Attribution

Yu Zheng* Boyang Gong* Fanye Kong Yueqi Duan Bingyao Yu Wenzhao Zheng Lei Chen†
Jiwen Lu Jie Zhou
Tsinghua University, China

{yu-zheng, duanyueqi, yuby, leichenth, lujiwen, jzhou}@tsinghua.edu.cn;
wenzhao.zheng@outlook.com; boyanggong429@gmail.com; kongfy23@mails.tsinghua.edu.cn

Abstract

In this paper, we propose a Counterfactually Decoupled Attention Learning (CDAL) method for open-world model attribution. Existing methods rely on handcrafted design of region partitioning or feature space, which could be confounded by the spurious statistical correlations and struggle with novel attacks in open-world scenarios. To address this, CDAL explicitly models the causal relationships between the attentional visual traces and source model attribution, and counterfactually decouples the discriminative model-specific artifacts from confounding source biases for comparison. In this way, the resulting causal effect provides a quantification on the quality of learned attention maps, thus encouraging the network to capture essential generation patterns that generalize to unseen source models by maximizing the effect. Extensive experiments on existing open-world model attribution benchmarks show that with minimal computational overhead, our method consistently improves state-of-the-art models by large margins, particularly for unseen novel attacks. Source code: <https://github.com/yzheng97/CDAL>.

1. Introduction

The rapid development of visual generative models [11, 12, 15, 22, 26, 28, 47, 50, 51, 55, 56] has raised increasing concerns about their potential misuse for personal reputation damage and public misinformation. Although deepfake detectors can effectively distinguish real from synthetic content [32, 35, 42, 54, 62, 69], such capability alone does not address the need to identify source generative models for further legal measures, referred to as model attribution.

Pioneering works on model attribution actively root fingerprints during training [27, 65–67] or parse the given AI-

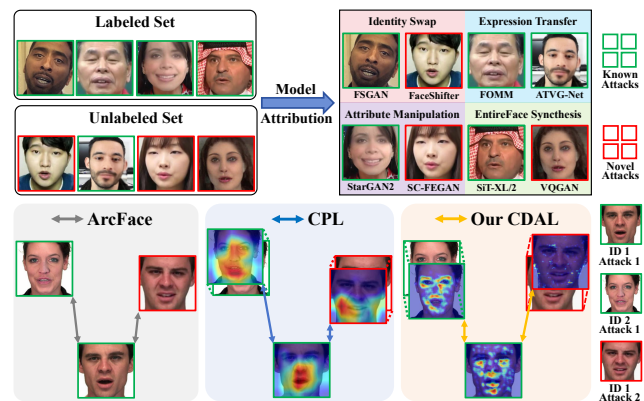


Figure 1. **Upper:** Problem Setup. **Bottom:** Exemplified sample in feature space. **Bottom left:** ArcFace [9] distances show that forgery images of the same source identity are clustered together regardless of their attacking models. **Bottom middle:** Feature differences learned by existing methods [52] are still influenced by source bias (i.e., identity) from the novel attack rather than model-specific artifacts. **Bottom right:** Our approach effectively isolates model-specific artifacts from source content bias, thus enabling accurate attribution even for unseen generative models.

generated images [4, 16, 30, 36, 39, 61, 63, 65]. However, their closed-set assumption becomes impractical as new generative models emerge rapidly and challenge their generalizability. To this end, researchers have built benchmarks on open-world model attribution and accordingly devised techniques aimed to attribute known attacks to their source models while identifying unseen novel attacks in real world scenarios (Figure 1 upper). For example, [52] designs a Contrastive Pseudo Learning framework for open-world deepfake attribution, which is later extended with multi-scale learning and frequency information [53]. [64] simulates a feature space spanning mechanism specifically designed for GANs. To exploit the subtle artifacts from various generative models, these approaches design handcrafted region partitions to search for useful local patterns [52, 53], or simulate the progressive feature space spanning via man-

*Equal contribution.

†Corresponding author.

ually modifying the model fingerprints [64].

However, as exemplified in Figure 1, apart from their inherent model-specific artifacts, generative models preserve semantic biases from source images during forgery processes. This poses a critical attribution challenge where the semantic content dominates over model-specific traces. The ArcFace [9] embeddings (Figure 1 bottom left) serve as evidence of this challenge where forgery images are clustered by identity rather than by generative source in feature space. Even for state-of-the-art methods like CPL [52] (Figure 1 bottom middle), despite their sophisticated design of handcrafted pixel partition or feature space, they remain confounded by source bias content. Their resulting attention maps fail to consistently highlight model-specific forgery patterns, thus leading to poor generalization when confronting novel unseen attacks where the previously learned spurious correlations no longer apply.

In this paper, we propose a plug-and-play Counterfactually Decoupled Attention Learning (CDAL) method for open-world model attribution. Different from previous methods confounded by the spurious correlations, CDAL aims to explicitly model the essential causal relationships between visual forgery traces and attributed source models that generalize well in open-world scenarios. Specifically, taking inspiration from counterfactual intervention, we isolate model-specific artifacts from source content biases by extracting factual and counterfactual attentions from the input feature representation. Given the diversified and uncertain nature of forgery traces in open-world scenarios, we further propose the Causal Attention Augmentation to achieve broader spatial coverage for these decoupled attentions while maintaining the constructed causal consistency. The quality of the learned attention maps can be quantified through the causal effect, i.e., the difference between the attribution results predicted by factual and counterfactual attentions. By maximizing this effect, CDAL provides a supervision signal to encourage the attribution network to attend to discriminative generation patterns in open-world model attribution (Figure 1 bottom right), while resisting misleading source biases. With negligible computational overhead, our method can be incorporated readily into existing methods where CDAL significantly improves baseline performances across different OWMA benchmarks, including OW-DFA [52] for deepfake attribution, and OSMA [64] for GAN attribution and discovery.

2. Related Works

Model Attribution: Model attribution aims to identify the source generative model of visual content. Active methods [27, 65–67] inject specifically-designed fingerprints during the training phase of the generative model. However, they cannot handle the “black-box” scenario without access to the source model and its training pipeline. In con-

trast, passive attribution methods [4, 30, 36, 61] focus on identifying the source model solely from its generated outputs. For example, [16, 39, 63, 65] attribute GAN architectures through parsing their inherent model fingerprints within the generated images. Our CDAL belongs to the model-agnostic passive technique, and is applicable to various scenarios including deepfake model attribution, GAN attribution and GAN discovery.

Open-world Recognition: Open-world Recognition aims at identifying known categories and simultaneously recognizing novel classes unseen during training, which has been explored in visual classification [5, 17, 46], semantic segmentation [49] and object detection [24], etc. As new generative models emerge rapidly, researches have paid attention to constructing benchmarks to simulate the open-world model attribution scenarios and developing the corresponding solutions [14, 52, 53, 64]. Dynamic network adaptation has shown promise where hypernetworks enable scene-conditioned learning to adapt models to different input contexts [70, 71]. However, these approaches primarily rely on statistical correlations through handcrafted designs of region partitioning [52, 53] or feature space [64], which limits their effectiveness when the learned correlations no longer hold for diverse unseen models. Our CDAL instead focuses on modeling causal relationships between visual forgery traces and source models, which enables more robust generalization to previously unseen forgery techniques.

Causal Reasoning in Computer Vision: There has been a number works leveraging causal reasoning [40, 43] to benefit various sub-tasks in computer vision, including classification [38], captioning [59], re-identification [33, 45], etc. To isolate the direct causal effect of interest critical to the target task, they perform counterfactual interventions through the *do*-operator, which is typically implemented by randomized attention maps to cut off the causal path between confounding variables and outcomes [33, 45, 58, 59]. We are inspired by this intervention operation and propose to decouple the factual and counterfactual attentions in the context of open-world model attribution.

3. Approach

In this section, we present our proposed Counterfactually Decoupled Attention Learning (CDAL) framework. We firstly introduce the addressed task from a causal perspective. Then we outline our CDAL with our core insights of counterfactually decoupled attention learning followed by the technical details of counterfactual feature isolation and causal attention augmentation.

3.1. Causal Perspective on the Addressed Problem

Problem Setup: Open-world model attribution [52, 53, 64] aims to identify the source generative model of synthetic images, including those from previously unseen models.

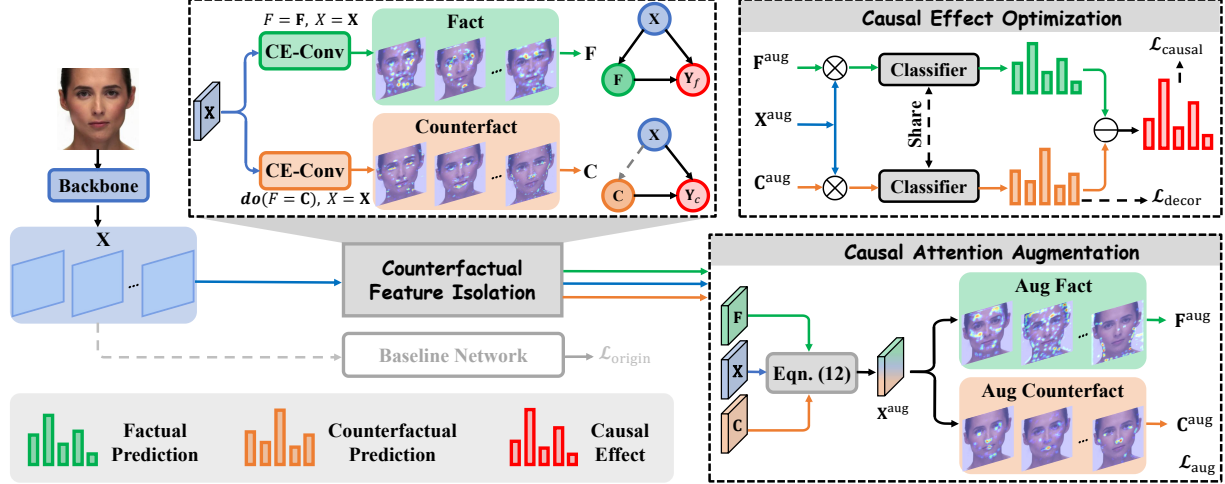


Figure 2. Overview of our proposed CDAL, which can be readily incorporated into existing baseline networks. **Upper right:** CDAL fundamentally maximizes the causal effect, i.e., the difference between the predictions from factual and counterfactual attentions, to encourage the network to learn more effective visual attention for model attribution and reduce the effects of biased training data. **Upper left:** The factual and counterfactual attentions are generated by Counterfactual Feature Isolation, which simulates counterfactual intervention and employs parallel Causal Expert (CE) convolutions. **Bottom right:** We complement the extracted attention maps with broader spatial coverage through Causal Attention Augmentation while maintaining the constructed causal consistency.

Given known generative models $\mathbb{G}_K = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_M\}$ and their generated images $\mathcal{X}_K = \{x_i | x_i \sim \mathcal{G}_j, \mathcal{G}_j \in \mathbb{G}_K\}$, we learn a function f that maps:

$$f : \mathcal{X} \rightarrow \mathbb{G}_K \cup \{\text{unknown}\} \quad (1)$$

When $f(x) = \text{unknown}$, the image is determined to be generated by some unknown model $\mathcal{G}_u \in \mathbb{G}_U$, where \mathbb{G}_U represents new generative models not available during training. As aforementioned, a key challenge lies in the problem that two types of features are blended in synthetic images:

- Model-specific artifacts that *causally relate to* source models.
- Source bias features from original content with *no causal link to* generation methods.

Existing approaches typically rely on handcrafted region partitioning [52, 53] or feature spanning [64] strategies, which could be confounded by the learned statistical correlations. As they fail to isolate the true causal relationships between visual artifacts and their generating processes, these spurious correlations might break down when confronted with previously unseen generative models.

Structural Causal Model for Attribution: To address this limitation, we reformulate the problem with a Structural Causal Model (SCM): $G = \{N, E\}$, where nodes N and edges E represent variables and their causal dependencies:

$$X \rightarrow A \rightarrow Y \text{ and } X \rightarrow Y. \quad (2)$$

Here, X, A, Y denotes the input image, feature maps used for attribution (specifically, the attention maps), and source model prediction. The attention map $\mathbf{A}_i \in \mathbb{R}_+^{H \times W}$ highlights regions in the image that are important for attribution decisions, where H and W are the height and width of the attention map. These attention maps are used to selectively

focus on discriminative regions by weighing feature maps through element-wise multiplication:

$$\mathbf{h}_i = \varphi(X * \mathbf{A}_i) = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W X^{h,w} \mathbf{A}_i^{h,w}, \quad (3)$$

where $*$ and φ denotes Hadamard product and a global average pooling operation that aggregates the weighted features. The attention-weighted representations from different regions are then combined to form the final representation for attribution prediction.

Our CDAL: Built upon the structural causal model above, our CDAL framework aims to capture the essential causal relationships between visual forgery traces and the attributed source models, whose overall framework is illustrated in Figure 2. The basic idea is to quantify the quality of the learned feature maps by comparing the effects of factual attention \mathbf{F} (i.e., the learned model-specific attention) and counterfactual attention \mathbf{C} (i.e., the attention on source bias). This difference (**upper-right** in Figure 2), also known as causal effect [40, 43], is maximized to encourage the network to learn more effective visual attention for model attribution and reduce the effects of biased training data.

Specifically, given the input feature \mathbf{X} , we are inspired by counterfactual intervention $do(A = \bar{A})$ [43] to isolate model-specific artifacts from source content bias (**upper left** in Figure 2), where we extract the factual attentions \mathbf{F} and counterfactual attentions \mathbf{C} from \mathbf{X} :

$$\begin{aligned} \mathbf{F} &= \mathcal{F}(\mathbf{X}) = \{\mathbf{F}_1, \dots, \mathbf{F}_M\}, \\ \mathbf{C} &= \mathcal{C}(\mathbf{X}) = \{\mathbf{C}_1, \dots, \mathbf{C}_M\}, \end{aligned} \quad (4)$$

where $\mathcal{F}(\cdot)$ and $\mathcal{C}(\cdot)$ extract factual and counterfactual attentions (**introduced in Section 3.2**) respectively. M is the

number of attention maps.

After obtaining these attention maps, given the uncertain and diverse nature of subtle forgery traces in open world, we complement them with broader spatial coverage while maintaining causal consistency (**bottom-right** in Figure 2). This is achieved through Causal Attention Augmentation (**introduced in Section 3.3**) which outputs the augmented features \mathbf{X}^{aug} and corresponding attention maps $\mathbf{F}^{\text{aug}}, \mathbf{C}^{\text{aug}}$. A shared attribution classifier $\delta(\cdot)$ is used to map them into factual and counterfactual predictions (\mathbf{Y}_f and \mathbf{Y}_c):

$$\mathbf{Y}_f = \mathbf{Y}(F=\mathbf{F}, X=\mathbf{X}) \triangleq \delta\left(\sum_{i=1}^M \mathbf{X}^{\text{aug}} * \mathbf{F}_i^{\text{aug}}\right), \quad (5)$$

$$\mathbf{Y}_c = \mathbf{Y}(\text{do}(F=\mathbf{C}), X=\mathbf{X}) \triangleq \delta\left(\sum_{i=1}^M \mathbf{X}^{\text{aug}} * \mathbf{C}_i^{\text{aug}}\right).$$

The difference $\mathbf{Y}_{\text{effect}} = \mathbf{Y}_f - \mathbf{Y}_c$, inspired by causal effect analysis [40, 43], quantifies the quality of learned feature representations by measuring how much the model-specific artifacts (factual) outperform the bias-prone features (counterfactual) in attribution predictions. We leverage this causal effect as the optimization target by employing a cross-entropy (CE) loss function $\mathcal{L}_{\text{causal}}$:

$$\mathcal{L}_{\text{causal}} = \text{CE}(\mathbf{Y}_{\text{effect}}, y), \quad (6)$$

where y is the ground-truth attribution label. This optimization target creates opposing learning objectives where factual and counterfactual attentions are driven toward discriminative model-specific traces and source-dependent biases respectively. By maximizing their performance gap, we strengthen the causal connection between the attentional features and source model attribution.

3.2. Counterfactual Feature Isolation

To extract the counterfactual attention maps in Eqn. (4), a common practice is to replace the original attention with the ones with randomized fixed weights [33, 58, 59]. However, such a predefined distribution on counterfactual attention weights, along with the handcrafted design [52, 53, 64] in search of factual forgery traces, still typically rely on statistical correlations that inadequately capture diverse artifact characteristics in unseen attacks (see Table 5d). To address this limitation, we employ Causal Expert Convolution that is guided by the aforementioned opposite learning objectives and actively learns to extract counterfactual patterns without predefined assumptions, rather than simply cut off the causal effect between input and counterfactual attention. **Causal Expert Convolution:** Unlike static convolutions that indiscriminately mix all feature correlations, we employ Causal Expert (CE) convolutions to explicitly model the causal effect through which model-specific artifacts influence the attribution prediction. This is achieved by dynamically constructing a convolution kernel W' that

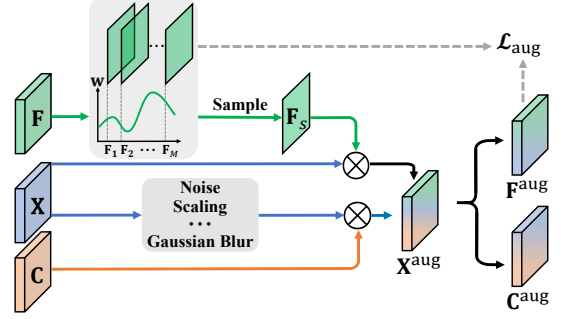


Figure 3. Illustration of Casual Attention Augmentation.

weights N_{exp} different expert kernels based on their causal relevance to the attribution task:

$$W' = \sum_{i=1}^{N_{exp}} \alpha_i W_i, \quad \alpha = \sigma(\text{MLP}(\text{Pool}(\mathbf{X}))) \in \mathbb{R}^{N_{exp}} \quad (7)$$

where α denotes the estimated causal contribution factors and W_i is the i -th expert kernel.

Causal Attention Generation: The CE convolution is used to extract both factual (F) and counterfactual (C) attention maps in Eqn. (4). Below we describe the process for generating F , which consists of two complementary extraction paths. First, a 1×1 CE convolution operator consumes \mathbf{X} to obtain $\mathbf{X}_{\text{cross}}$, which preserves global relationships across channels. $\mathbf{X}_{\text{cross}}$ is then fed into a depth-wise separable [7] CE convolution, producing $\mathbf{X}_{\text{depth}}$ with its channel-specific patterns. This design helps isolate channel-specific patterns [19] that may contain distinctive model artifacts, which might otherwise be diluted in standard convolutions. The final factual attention map combines these complementary features:

$$\mathbf{F} = \text{Concat}[\mathbf{X}_{\text{cross}}, \mathbf{X}_{\text{depth}}]. \quad (8)$$

Similarly, the counterfactual maps \mathbf{C} are generated by this two-stage process with another set of network parameters.

Since forgery methods preserve semantic source bias while introducing model-specific artifacts, we enhance the separation between discriminative features and source biases through a decorrelation loss:

$$\mathcal{L}_{\text{decor}} = \text{CE}(\mathbf{Y}_c, \mathbf{Y}_c) = - \sum_{c \in \mathbf{C}} \mathbf{Y}_c \log \mathbf{Y}_c. \quad (9)$$

By maximizing the entropy in counterfactual predictions, $\mathcal{L}_{\text{decor}}$ pushes \mathbf{Y}_c towards a uniform distribution across attributed classes, and forces counterfactual attention to focus on non-causal source content that lacks discriminative value for model attribution.

3.3. Causal Attention Augmentation

While our attention learning aims to help identify potential forgery traces, two critical challenges remain in open-world scenarios: (1) the uncertainty in unseen source features (new faces, unexpected noises) that can weaken decoupling efficacy, and (2) the inherent diversity of forgery

traces in full-face synthesis which can appear across multiple locations. To address these challenges, we propose a Causal Attention Augmentation operation (shown in Figure 3) that expands attention coverage while maintaining causal relationships. It firstly generates diversified feature representations through standard augmentation techniques:

$$\mathbf{X}^{\text{aug}} = \mathcal{A}_1(\mathcal{A}_2(\cdots \mathcal{A}_N(\mathbf{X}))), \quad (10)$$

where $\mathcal{A}_1, \cdots, \mathcal{A}_N$ are a series of operations such as noise adding or Gaussian blurring (see Appendix for details).

However, indiscriminative application of these augmentations risks diluting model-specific traces and corrupting established causal relationships. We therefore employ a targeted augmentation strategy to preserve forgery-relevant regions while diversifying source-specific features only. To explore diverse forgery traces while maintaining causal consistency, we probabilistically select an influential factual attention map indexed by s :

$$s = \arg, \underset{i \in \{1, 2, \dots, M\}}{\text{sample}}(i \mid p(i) = \mathbf{w}_i) \quad (11)$$

where $\mathbf{w} \in \mathbb{R}^M$ is the normalized energy distribution of factual attention channels (see Appendix for derivation).

The corresponding sampled attention map \mathbf{F}_s is then used for selective augmentation:

$$\mathbf{X}^{\text{aug}} \leftarrow \mathbf{X} * \mathbf{F}_s + \mathbf{X}^{\text{aug}} * \mathbf{C}, \quad (12)$$

which aims to maintain the consistency of augmented samples with original samples in factual regions, while allowing counterfactual regions to exhibit diverse distributions. Then we can re-generate the enhanced factual attention $\mathbf{F}^{\text{aug}} = \mathcal{F}(\mathbf{X}^{\text{aug}})$ and counterfactual attention $\mathbf{C}^{\text{aug}} = \mathcal{C}(\mathbf{X}^{\text{aug}})$ by re-using the networks in Section 3.2.

To prevent the model from over-concentrating on single attention areas and to encourage the exploration of alternative forgery manifestations, we further introduce an attention diversification loss \mathcal{L}_{aug} as follows:

$$\mathcal{L}_{\text{aug}} = \frac{1}{M} \sum_{i=1}^M |\mathbf{X} * \mathbf{F}_i - \mathbf{X}^{\text{aug}} * \mathbf{F}_i^{\text{aug}}| \cdot (1 - \mathbf{1}_{i=s}), \quad (13)$$

where $\mathbf{1}_{(\cdot)}$ is the indicator function. M is the number of factual attention maps. Through minimizing \mathcal{L}_{aug} , it maintains causal consistency [68] between original and augmented features in model-specific regions while encouraging exploration of complementary forgery traces by excluding the already-attended regions from the consistency constraint.

The final loss function $\mathcal{L}_{\text{total}}$ is computed as follows:

$$\mathcal{L}_{\text{total}} = \eta_1 \mathcal{L}_{\text{causal}} + \eta_2 \mathcal{L}_{\text{decor}} + \eta_3 \mathcal{L}_{\text{aug}} + \mathcal{L}_{\text{original}} \quad (14)$$

where η_1 , η_2 and η_3 are hyper-parameters. $\mathcal{L}_{\text{original}}$ is the original loss in baselines which also takes \mathbf{X} as input.

4. Experiments

In this section, we evaluate CDAL on the public benchmarks for open-world model attribution, including open-

world deepfake attribution [52], GAN attribution and discovery [64]. Unlike closed-set attribution where source models of test samples are present in the training set, these benchmarks challenge the model to generalize to previously unseen generative models. The key challenge lies in preventing the learned features or patterns from overfitting to known models while maintaining the discrimination power. We incorporated CDAL into state-of-the-art approaches in each benchmark (see Appendix for details), and conducted comprehensive comparisons with various baseline methods. In addition, we performed extensive ablation studies to analyze the effect of key components, hyper-parameters, design choices, as well as the model efficiency of our method.

4.1. Open-world Deepfake Attribution

Dataset: The task aims to simultaneously attribute fake face images to known source models and identify those from unknown ones. We experimented on the OW-DFA dataset [52] which contains 20 challenging forgery methods and real face images from FaceForensics++ [48] and Celeb-DF [34], whose forgery types include identity swap, expression transfer, attribute manipulation, entire face synthesis. Each type (including real faces) contains both labeled and unlabeled images.

Protocols: We followed [52] to employ two protocols. Protocol 1 focuses on forged images only, where labeled and unlabeled images were treated as known and novel attacks encountered in the open world. In Protocol 2, real images from FaceForensics++ [48] (labeled) and Celeb-DF [34] (unlabeled) are blended with forged images to create a dataset where real faces substantially outnumber fake ones as in real-world conditions.

Evaluation Metrics: We followed [52] to use classification Accuracy (ACC), Normalized Mutual Information (NMI), and Adjusted Rand Index (ARI) to evaluate on the OW-DFA dataset, where Hungarian algorithm [29] was applied to align predictions with ground-truth labels.

Results: To get deeper insights into CDAL, we firstly visualize the learned attentions via CAM [72] in Figure 4. While CPL [52] attends to broader, less focused regions that potentially dilute discriminative traces, our factual attention identifies model-specific patterns critical for attribution. Meanwhile, the counterfactual attention highlights source bias regions (including irrelevant background) that would otherwise introduce spurious correlations and confound the attribution process.

The quantitative results on OW-DFA are shown in Table 1, where the ‘‘upper bound’’ and ‘‘lower bound’’ represent supervised learning on the labeled set and labeled+unlabeled sets respectively [52]. In Protocol 1, our method demonstrates significant advantages in open-world scenarios. Particularly, CDAL significantly boosts the novel attack attribution by 11.27% ARI for CPL [52], which out-

Table 1. Quantitative results (%) of **Protocol 1** and **Protocol 2** on OW-DFA [52].

Method	Protocol-1: Fake						Protocol-2: Real & Fake							
	Known		Novel		All		Known		Novel		All			
	ACC	ACC	NMI	ARI	ACC	NMI	ARI	ACC	ACC	NMI	ARI	ACC	NMI	ARI
Lower Bound	99.68	40.86	47.55	26.33	46.91	63.43	37.33	99.84	34.57	42.98	19.37	61.46	66.05	62.16
Upper Bound	98.93	96.99	94.18	94.94	97.91	95.87	95.91	99.27	97.12	94.89	96.78	98.43	96.48	98.27
Openworld-GAN [14]	99.49	39.14	47.08	44.39	57.92	58.24	48.71	99.51	40.09	50.72	35.96	67.02	57.81	59.92
DNA-Det [63]	99.62	38.16	48.76	23.21	46.88	66.15	35.46	99.45	39.03	47.07	22.54	60.68	67.94	56.49
RankStats [18]	99.17	62.05	64.60	52.87	79.52	78.87	72.90	98.86	51.19	57.56	37.56	78.25	77.37	88.07
OpenLDN [46]	98.78	54.12	57.54	45.43	72.90	77.22	70.03	97.03	48.26	52.77	33.72	73.97	75.13	84.37
ORCA [5]	98.30	73.61	70.20	63.50	85.23	83.99	80.86	97.09	62.10	64.96	49.15	83.44	82.68	88.64
NACH [17]	98.34	73.43	71.61	65.33	85.16	84.90	82.31	97.28	69.39	70.03	54.28	86.47	84.76	90.09
CPL [52]	98.68	75.21	73.19	65.71	86.25	85.58	82.35	97.45	69.57	70.67	54.67	86.51	85.44	90.30
MPSL [53]	98.55	75.23	76.99	69.26	86.31	87.27	84.91	98.01	69.81	72.77	55.97	86.71	86.60	92.37
ORCA + Ours	98.55	76.85	77.18	67.77	86.91	88.04	83.74	98.44	66.05	70.65	50.42	86.20	86.59	91.65
<i>Improvement</i>	+0.25	+3.24	+6.98	+4.27	+1.68	+4.05	+2.88	+1.35	+3.95	+5.69	+1.27	+2.76	+3.91	+3.01
NACH + Ours	98.68	85.31	81.22	74.50	91.60	90.25	86.74	98.79	69.92	73.93	56.90	87.44	88.04	92.28
<i>Improvement</i>	+0.34	+11.88	+9.61	+9.17	+6.44	+5.35	+4.43	+1.51	+0.53	+3.90	+2.62	+0.97	+3.28	+2.19
CPL + Ours	98.90	86.02	82.19	76.98	92.06	90.60	87.66	99.08	70.56	74.14	57.22	88.02	88.59	93.01
<i>Improvement</i>	+0.22	+10.81	+9.00	+11.27	+5.81	+5.02	+5.31	+1.63	+0.99	+3.47	+2.55	+1.51	+3.15	+2.71

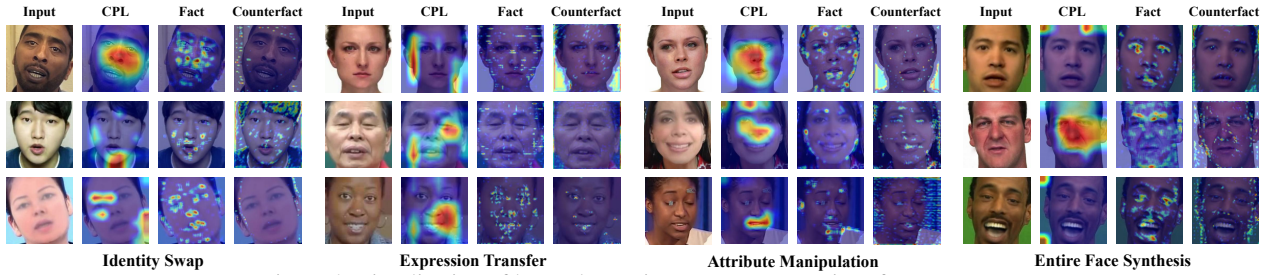


Figure 4. Visualization of learned attention maps across various forgery types.

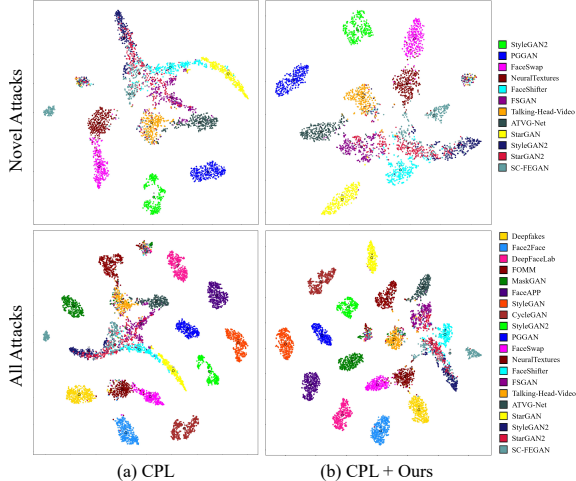


Figure 5. t-sne visualization in OW-DFA.

performs previous state-of-the-art MPSL [53] by 7.72%. The improvements upon ORCA and NACH, while not designed for model attribution, collectively validate the generalization ability of CDAL. In Protocol-2 which involves real faces simulating challenging real-world scenarios, our method also contributes to strong performance gains. Compared to MPSL [53], CPL [52]+Ours leads by 1.37% in NMI for novel attacks. Integration into other baselines like ORCA [5] also improves the key metrics, which validates

Table 2. Results (%) of **Protocol 1** on OW-DFA [52] extended with modern generative models from DF40 [62].

Setting	Method	Known			Novel			All		
		ACC	ACC	NMI	ARI	ACC	NMI	ARI		
Setting1	CPL [52]	98.83	76.40	78.16	65.96	82.94	87.50	78.76		
	CPL + Ours	99.80	83.07	82.73	74.17	86.75	89.75	82.52		
	<i>Improvement</i>	+0.97	+6.67	+4.57	+8.21	+3.81	+2.25	+3.76		
Setting2	CPL [52]	99.21	76.31	74.04	66.45	87.10	86.87	83.28		
	CPL + Ours	99.34	80.00	79.53	72.64	89.07	89.30	85.90		
	<i>Improvement</i>	+0.13	+3.69	+5.49	+6.19	+1.97	+2.43	+2.62		

our proposed method in real-world scenarios.

t-SNE Visualization: To visually compare CDAL and baseline methods, we presented t-SNE [57] results in Figure 5 using the features for final prediction. The visualization includes 8 known attacks and 12 novel attacks. CDAL significantly improves the clustering performance on novel attacks while maintaining discriminative power for known attacks. For novel attacks [6, 23, 31, 41], CDAL exhibits effective discrimination ability. Even in scenarios involving cross-source datasets (e.g., DeepFakes [2] from FaceForensics++ [48] and DeepFaceLab [1] from ForgeryNet [20] both belonging to identity swap; PGGAN [25] from DFFD [8] and CycleGAN [73] from ForgeryNIR [60] both belonging to Entire Face Synthesis), CDAL still substantially reduces the feature distances of the samples from the same forgery type. See Appendix for the t-SNE visualizations under Protocol 2.

Table 3. Quantitative results (%) of GAN attribution (Protocol 1) on OSMA [64], which are averaged among five splits.

Method	Closed-Set	Unseen Seed		Unseen Architecture		Unseen Dataset		Unseen All	
	Accuracy	AUC	OSCR	AUC	OSCR	AUC	OSCR	AUC	OSCR
PRNU [39]	55.27	69.20	49.16	70.02	49.49	67.68	48.57	68.94	49.06
Yu et al [65]	85.71	53.14	50.99	69.04	64.17	78.79	72.20	69.90	64.86
DCT-CNN [13]	86.16	55.46	52.68	72.56	67.43	72.87	67.57	69.46	64.70
DNA-Det [63]	93.56	61.46	59.34	80.93	76.45	66.14	63.27	71.40	68.00
RepMix [4]	93.69	54.70	53.26	72.86	70.49	78.69	76.02	71.74	69.43
POSE [64]	94.81	68.15	67.25	84.17	81.62	88.24	85.64	82.76	80.50
RepMix + Ours	94.01	59.20	57.41	75.38	73.11	81.73	78.96	73.97	71.72
<i>Improvement</i>	+0.32	+4.50	+4.15	+2.52	+2.62	+3.04	+2.94	+2.23	+2.29
POSE + Ours	95.25	72.99	71.73	86.97	84.48	91.32	88.66	85.37	82.85
<i>Improvement</i>	+0.44	+4.84	+4.48	+2.80	+2.86	+3.08	+3.02	+2.61	+2.35

Generalization to Modern Generative Models: To keep pace with the development tendency, we extended the OW-DFA benchmark with diffusion-based and flow-based models (see Appendix for details). We design two experimental settings: In Setting 1, we augment both training and test sets with diffusion models [3, 37, 44] and flow-based model [12]) from the recent large-scale DF-40 dataset [62], where DiT-XL/2 [44] serves as a labeled known attack and the others as unlabeled novel attacks. As shown in Table 2, our method achieved substantial performance gains on novel attacks. In Setting 2, we evaluate the generalization on the completely unseen DDPM model [21]. Our method demonstrates superior performance compared to CPL [52]. This is because the CPL model learns the distribution within the training set, thus struggling with out-of-distribution attacks in this open-set scenario. These results validate our approach at both identifying known manipulation sources and discovering emerging generation techniques in real-world applications.

4.2. Open-world GAN Attribution

Dataset: This task aims to simultaneously attribute GAN-generated images to known GAN models and discover novel GAN classes in an open-world scenario. We experimented on the OSMA [64] benchmark with five splits, which contains 15 known classes as the close-set (including real and 14 seen models), and 53 unknown classes as the open-set in total. The 53 unknown classes comprise of models trained with 10 random seeds, 22 architecture and 21 dataset distributions respectively.

Protocols: We followed the protocols from [64]. In Protocol 1 (GAN Attribution), models were trained on the closed-set, and attribution confidence scores were computed for both closed-set and open-set using standard OSCR. In Protocol 2 (GAN Discovery), models were trained on known GAN classes. During testing, features of closed-set and open-set images were extracted and clustered via K-Means. Closed-set samples were mapped to known GAN classes, while open-set samples formed auxiliary clusters represented as novel classes.

Table 4. Quantitative results (%) of GAN discovery (Protocol 2) on OSMA [64], which are averaged among five splits.

Method	Close-set	Unseen All		
	ACC	Purity	NMI	ARI
RepMix [4]	94.01	31.53	51.60	18.71
POSE [64]	94.81	41.04	60.59	26.39
RepMix + Ours	93.83	37.96	52.08	20.66
<i>Improvement</i>	+0.32	+6.43	+0.48	+1.95
POSE + Ours	95.25	48.93	61.89	29.65
<i>Improvement</i>	+0.44	+7.89	+1.30	+3.26

Evaluation Metrics: Following [14, 64], for Protocol 1, we reported Area Under the ROC Curve (AUC) and Open Set Classification Rate (OSCR) [10], which balances the accurate classification of known models, and the correct discrimination between known and unknown models. For Protocol 2, we evaluated closed-set performance using Accuracy (ACC) and open-set performance using Average Purity, NMI, and ARI. Results were averaged on five splits [64].

Results: Results of Protocol 1 on OSMA are shown in Table 3. CDAL achieves the greatest improvements on the unseen seed setting, which indicates that CDAL can still effectively capture subtle differences due to randomness. Table 4 presents the comparison on Protocol 2 of GAN discovery. Our method improves significantly on both RepMix [4] and POSE [64], with overall purity in the unseen setting increased by 6.43 and 7.89 % respectively. These results highlight the effectiveness of CDAL in improving generalization to unseen GAN models across various scenarios. For qualitative results, please refer to the Appendix pages for the t-SNE visualization on the comparison of with and without our CDAL.

4.3. Ablation Studies and Analysis

In this subsection, we present ablation studies to verify our key components, hyper-parameters and model efficiency. See Appendix for further results and analysis.

Ablation Study on Key Components: Results on ablating the key components in CDAL is shown in Table 5a. Here, FA, CA, EA stand for factual attentions, counterfactual attentions, and the enhanced attentions after being

Table 5. Results of ablation studies and in-depth analysis of CDAL.

(a) Ablation study on key components of CDAL.								(b) Ablation study on loss functions of CDAL.								(c) Ablation study on N.							
Baseline	FA	CA	EA	Novel (OW-DFA)		Unseen All (OSMA)		Baseline	$\mathcal{L}_{\text{causal}}$	$\mathcal{L}_{\text{decor}}$	\mathcal{L}_{aug}	Novel (OW-DFA)		Unseen All (OSMA)		N	Known		Novel				
				ACC	NMI	ARI	Purity					NMI	ARI	ACC	NMI		ARI	Purity	NMI	ARI	ACC	ACC	NMI
✓				75.21	73.19	65.71	41.04	60.59	26.39				79.56	77.24	68.94	43.67	60.86	27.08	2	98.76	84.27	79.23	72.57
✓	✓			81.78	78.20	70.34	45.58	60.96	27.63				82.46	79.90	71.75	46.23	61.02	27.91	3	98.25	82.10	80.52	72.60
✓	✓	✓		84.05	80.86	73.54	46.65	61.07	27.95				84.85	81.80	75.08	47.22	61.34	28.19	4	98.90	86.02	82.19	76.98
✓	✓	✓	✓	86.02	82.19	76.98	48.93	61.89	29.65				86.02	82.19	76.98	48.93	61.89	29.65	5	98.68	82.68	81.94	74.56
																			6	98.68	79.77	79.94	72.00

(d) Comparison on counterfactual attentions.					(e) Comparison with vanilla attention.					(f) Comparison on model efficiency.				
Method	Known		Novel		Method	Known		Novel		OW-DFA				
	ACC	ACC	NMI	ARI		ACC	ACC	NMI	ARI	Method	ARI	Params/M	FLOPs/G	
CPL [52]	98.68	75.21	73.19	65.71	CPL [52]	98.68	75.21	73.19	65.71	CPL [52]	65.71	23.59	5.397	
CPL + Ours (Random)	98.35	84.05	80.86	73.54	CPL + Attention	98.25	70.37	70.59	58.81	CPL + Ours	76.98	23.91	5.399	
CPL + Ours (Uniform)	98.68	82.68	79.52	70.81	Δ	-0.43	-4.84	-2.60	-6.90	OSMA				
CPL + Ours (Reversed)	98.46	82.30	79.27	71.88	CPL + Ours	98.90	86.02	82.19	76.98	POSE [64]	26.39	22.68	1.039	
CPL + Ours (Shuffle)	98.46	80.74	79.56	72.12	Δ	+0.22	+10.81	+9.00	+11.27	POSE + Ours	29.65	22.93	1.041	
CPL + Ours (CE-Conv)	98.90	86.02	82.19	76.98										

augmented. When only FA is used, the model dynamically extracts the factual attention, which significantly improves performance on OW-DFA and OSMA. However, the limited improvement for unseen classes is attributed to the static random counterfactual attention. After introducing CA, the NMI for novel classes in OW-DFA and unseen classes in OSMA increases by 2.66% and 2.86%, respectively. This indicates a stronger adaptability to unseen classes. Finally, the introduction of EA achieves the best overall performance. EA, based on causal consistency, integrates factual and counterfactual attention to generate more robust feature representations. These results collectively validate the effectiveness of the designed component in CDAL.

Ablation Study on Loss Functions of CDAL: Table 5b shows the ablation results on loss functions. Compared to baseline with $\mathcal{L}_{\text{original}}$ only, $\mathcal{L}_{\text{causal}}$ achieves an improvement of +2.90% ACC on OW-DFA by decoupling factual attentions from counterfactual ones, which validates our overall causal modeling. This is further enhanced by $\mathcal{L}_{\text{decor}}$ which enables counterfactual attention to focus on source biases that might be shared across attribution models. \mathcal{L}_{aug} further masks the already attended regions to encourage the model to explore complementary feature regions, which leads to a notable gain of +1.9% ARI on OW-DFA. These ablational experiments collectively demonstrate the efficacy of the designed loss functions in our CDAL.

Ablation Study on the Number of Experts: Table 5c shows the results on OW-DFA Protocol-1 with varying numbers of experts (N) in CE Convolution. CDAL achieves the best performances across all metrics when $N=4$, which outperforms all the other trials by sizable margins.

Ablation Study on Type of Counterfactual Attention: Table 5d presents the results of different counterfactual attention types (see Appendix for definitions). Our strategy significantly outperforms both the baseline and static alternatives like random and uniform attention. This is because CE-Conv actively learns to extract adaptive counterfactual patterns, rather than simply applying fixed interven-

tions that cannot effectively address the diverse characteristics of different forgery methods in open-world scenarios.

Can Vanilla Attention Promote Attribution Performances? It is crucial to validate that it is our proposed design in CDAL that indeed promotes the performances, rather than the vanilla attention itself. We accordingly conducted experiments by removing CDAL and only adding the vanilla attention to baselines. From Table 5e, we can see that adding vanilla attention to CPL leads to performance declines across all metrics, as apposed to the significant improvements by CDAL. This degradation can be attributed to vanilla attention focusing excessively on sample-specific features, which introduces bias and weakends generalization. On the contrary, CDAL effectively filters out these biases by focusing on model-specific artifacts, leading to more robust and generalizable attribution ability.

Computational Overhead Analysis: Table 5f shows the computational overhead of our CDAL when incorporated into the baselines. With significant performance improvements, our method only brings with up to 0.35M additional network parameters and 0.002 GFLOPs, which demonstrates the high efficiency of our method.

5. Conclusion

In this paper, we propose Counterfactually Decoupled Attention Learning (CDAL) for open-world model attribution. Unlike existing methods that rely on handcrafted strategies susceptible to spurious correlations, CDAL explicitly models causal relationships between visual forgery traces and source models, which effectively decouples model-specific artifacts from source biases. By maximizing the causal effect between factual and counterfactual attention maps, our approach encourages networks to capture generalizable generation patterns. Experiments show CDAL consistently improves state-of-the-art models with minimal overhead, especially for unseen attacks. Future work could explore extending this framework to broader forensics tasks such as video and multi-modality attributions.

Acknowledgments

This work was supported in part by the National Natural Science Foundation of China under Grant 62125603, Grant 62306031, Grant 62336004, Grant 62321005, and Grant 62441616, and in part by the China Postdoctoral Science Foundation under Grant 2024M761674.

References

- [1] Deepfacelab. <https://github.com/iperov/DeepFaceLab>. Accessed: 2023-2-28. 6
- [2] Deepfakes. <https://github.com/deepfakes/faceswap>. Accessed: 2023-2-28. 6
- [3] Sara Atito, Muhammad Awais, and Josef Kittler. Sit: Self-supervised vision transformer. *arXiv preprint arXiv:2104.03602*, 2021. 7
- [4] Tu Bui, Ning Yu, and John Collomosse. Repmix: Representation mixing for robust attribution of synthesized images. In *ECCV*, pages 146–163, 2022. 1, 2, 7
- [5] Kaidi Cao, Maria Brbic, and Jure Leskovec. Open-world semi-supervised learning. In *ICLR*, 2022. 2, 6
- [6] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, pages 8789–8797, 2018. 6
- [7] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, pages 1251–1258, 2017. 4
- [8] Hao Dang, Feng Liu, Joel Stehouwer, Xiaoming Liu, and Anil K Jain. On the detection of digital face manipulation. In *CVPR*, pages 5781–5790, 2020. 6
- [9] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, pages 4690–4699, 2019. 1, 2
- [10] Akshay Raj Dhamija, Manuel Günther, and Terrance Boult. Reducing network agnostophobia. In *NeurIPS*, pages 9175–9186, 2018. 7
- [11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 34:8780–8794, 2021. 1
- [12] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, pages 12873–12883, 2021. 1, 7
- [13] Joel Frank, Thorsten Eisenhofer, Lea Schönherr, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. Leveraging frequency analysis for deep fake image recognition. In *ICML*, pages 3247–3258, 2020. 7
- [14] Sharath Girish, Saksham Suri, Sai Saketh Rambhatla, and Abhinav Shrivastava. Towards discovery and attribution of open-world gan generated images. In *ICCV*, pages 14094–14103, 2021. 2, 6, 7
- [15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 27, 2014. 1
- [16] Luca Guarnera, Oliver Giudice, Matthias Nießner, and Sebastiano Battiato. On the exploitation of deepfake model recognition. In *CVPRW*, pages 61–70, 2022. 1, 2
- [17] Lan-Zhe Guo, Yi-Ge Zhang, Zhi-Fan Wu, Jie-Jing Shao, and Yu-Feng Li. Robust semi-supervised learning when not all classes have labels. In *NeurIPS*, 2022. 2, 6
- [18] Kai Han, Sylvestre-Alvise Rebuffi, Sebastien Ehrhardt, Andrea Vedaldi, and Andrew Zisserman. Automatically discovering and learning new visual categories with ranking statistics. In *ICLR*, 2020. 6
- [19] Kai Han, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. Ghostnet: More features from cheap operations. In *CVPR*, pages 1580–1589, 2020. 4
- [20] Yinan He, Bei Gan, Siyu Chen, Yichun Zhou, Guojun Yin, Luchuan Song, Lu Sheng, Jing Shao, and Ziwei Liu. Forgerynet: A versatile benchmark for comprehensive forgery analysis. In *CVPR*, pages 4360–4369, 2021. 6
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020. 7
- [22] Zhizhong Huang, Shouzheng Chen, Junping Zhang, and Hongming Shan. Pfa-gan: Progressive face aging with generative adversarial network. *TIFS*, 16:2031–2045, 2020. 1
- [23] Youngjoo Jo and Jongyoul Park. Sc-fegan: Face editing generative adversarial network with user’s sketch and color. In *ICCV*, pages 1745–1753, 2019. 6
- [24] KJ Joseph, Salman Khan, Fahad Shahbaz Khan, and Vineeth N Balasubramanian. Towards open world object detection. In *CVPR*, pages 5830–5840, 2021. 2
- [25] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018. 6
- [26] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410, 2019. 1
- [27] Changhoon Kim, Yi Ren, and Yezhou Yang. Decentralized attribution of generative models. In *ICLR*, 2021. 1, 2
- [28] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 1
- [29] HW Kuhn et al. The hungarian method for the assignment problem. *NRL*, 2(1-2):83–97, 1955. 5
- [30] Jialiang Li, Haoyue Wang, Sheng Li, Zhenxing Qian, Xinpeng Zhang, and Athanasios V Vasilakos. Are handcrafted filters helpful for attributing ai-generated images? In *ACM MM*, pages 10698–10706, 2024. 1, 2
- [31] Lingzhi Li, Jianmin Bao, Hao Yang, Dong Chen, and Fang Wen. Faceshifter: Towards high fidelity and occlusion aware face swapping. *arXiv preprint arXiv:1912.13457*, 2019. 6
- [32] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-ray for more general face forgery detection. In *CVPR*, pages 5001–5010, 2020. 1
- [33] Xulin Li, Yan Lu, Bin Liu, Yating Liu, Guojun Yin, Qi Chu, Jinyang Huang, Feng Zhu, Rui Zhao, and Nenghai

- Yu. Counterfactual intervention feature transfer for visible-infrared person re-identification. In *ECCV*, pages 381–398, 2022. 2, 4
- [34] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics. In *CVPR*, pages 3207–3216, 2020. 5
- [35] Bo Liu, Fan Yang, Xiuli Bi, Bin Xiao, Weisheng Li, and Xinbo Gao. Detecting generated images by real images. In *ECCV*, pages 95–110, 2022. 1
- [36] Fengyuan Liu, Haochen Luo, Yiming Li, Philip Torr, and Jindong Gu. Which model generated this image? a model-agnostic approach for origin attribution. In *ECCV*, pages 282–301, 2024. 1, 2
- [37] Jiawei Liu, Qiang Wang, Huijie Fan, Yinong Wang, Yandong Tang, and Liangqiong Qu. Residual denoising diffusion models. In *CVPR*, pages 2773–2783, 2024. 7
- [38] David Lopez-Paz, Robert Nishihara, Soumith Chintala, Bernhard Scholkopf, and Léon Bottou. Discovering causal signals in images. In *CVPR*, pages 6979–6987, 2017. 2
- [39] Francesco Marra, Diego Gragnaniello, Luisa Verdoliva, and Giovanni Poggi. Do gans leave artificial fingerprints? In *MIPR*, pages 506–511, 2019. 1, 2, 7
- [40] Leland Gerson Neuberg. Causality: models, reasoning, and inference, by judea pearl, cambridge university press, 2000. *Econometric Theory*, 19(4):675–685, 2003. 2, 3, 4
- [41] Yuval Nirkin, Yosi Keller, and Tal Hassner. FSGAN: Subject agnostic face swapping and reenactment. In *ICCV*, pages 7184–7193, 2019. 6
- [42] Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. Towards universal fake image detectors that generalize across generative models. In *CVPR*, pages 24480–24489, 2023. 1
- [43] Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic books, 2018. 2, 3, 4
- [44] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, pages 4195–4205, 2023. 7
- [45] Yongming Rao, Guangyi Chen, Jiwen Lu, and Jie Zhou. Counterfactual attention learning for fine-grained visual categorization and re-identification. In *ICCV*, pages 1025–1034, 2021. 2
- [46] Mamshad Nayeem Rizve, Navid Kardan, Salman Khan, Fahad Shahbaz Khan, and Mubarak Shah. Openldn: Learning to discover novel classes for open-world semi-supervised learning. In *ECCV*, pages 382–401, 2022. 2, 6
- [47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 1
- [48] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *ICCV*, pages 1–11, 2019. 5, 6
- [49] Matteo Sodano, Federico Magistri, Lucas Nunes, Jens Behley, and Cyrill Stachniss. Open-world semantic segmentation including class similarity. In *CVPR*, pages 3184–3194, 2024. 2
- [50] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *NeurIPS*, 28, 2015. 1
- [51] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 1
- [52] Zhimin Sun, Shen Chen, Taiping Yao, Bangjie Yin, Ran Yi, Shouhong Ding, and Lizhuang Ma. Contrastive pseudo learning for open-world deepfake attribution. In *ICCV*, pages 20882–20892, 2023. 1, 2, 3, 4, 5, 6, 7, 8
- [53] Zhimin Sun, Shen Chen, Taiping Yao, Ran Yi, Shouhong Ding, and Lizhuang Ma. Rethinking open-world deepfake attribution with multi-perspective sensory learning. *IJCV*, pages 1–24, 2024. 1, 2, 3, 4, 6
- [54] Chuangchuan Tan, Yao Zhao, Shikui Wei, Guanghua Gu, Ping Liu, and Yunchao Wei. Rethinking the up-sampling operations in cnn-based generative network for generalizable deepfake detection. In *CVPR*, pages 28130–28139, 2024. 1
- [55] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *NeurIPS*, 37:84839–84865, 2024. 1
- [56] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *NeurIPS*, 30, 2017. 1
- [57] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 9(11), 2008. 6
- [58] Hanqing Wang, Wei Liang, Jianbing Shen, Luc Van Gool, and Wenguan Wang. Counterfactual cycle-consistent learning for instruction following and generation in vision-language navigation. In *CVPR*, pages 15471–15481, 2022. 2, 4
- [59] Tan Wang, Jianqiang Huang, Hanwang Zhang, and Qianru Sun. Visual commonsense r-cnn. In *CVPR*, pages 10760–10770, 2020. 2, 4
- [60] Yukai Wang, Chunlei Peng, Decheng Liu, Nannan Wang, and Xinbo Gao. Forgeryinir: deep face forgery and detection in near-infrared scenario. *TIFS*, 17:500–515, 2022. 6
- [61] Zhenting Wang, Chen Chen, Yi Zeng, Lingjuan Lyu, and Shiqing Ma. Where did i come from? origin attribution of ai-generated images. *NeurIPS*, 36:74478–74500, 2023. 1, 2
- [62] Zhiyuan Yan, Taiping Yao, Shen Chen, Yandan Zhao, Xinghe Fu, Junwei Zhu, Donghao Luo, Chengjie Wang, Shouhong Ding, Yunsheng Wu, et al. Df40: Toward next-generation deepfake detection. In *NeurIPS*, 2024. 1, 6, 7
- [63] Tianyun Yang, Ziyao Huang, Juan Cao, Lei Li, and Xirong Li. Deepfake network architecture attribution. In *AAAI*, pages 4662–4670, 2022. 1, 2, 6, 7
- [64] Tianyun Yang, Danding Wang, Fan Tang, Xinying Zhao, Juan Cao, and Sheng Tang. Progressive open space expansion for open-set model attribution. In *CVPR*, pages 15856–15865, 2023. 1, 2, 3, 4, 5, 7, 8
- [65] Ning Yu, Larry S Davis, and Mario Fritz. Attributing fake images to gans: Learning and analyzing gan fingerprints. In *ICCV*, pages 7556–7566, 2019. 1, 2, 7
- [66] Ning Yu, Vladislav Skripniuk, Sahar Abdelnabi, and Mario Fritz. Artificial fingerprinting for generative models: Rooting deepfake attribution in training data. In *ICCV*, pages 14448–14457, 2021.
- [67] Ning Yu, Vladislav Skripniuk, Dingfan Chen, Larry Davis, and Mario Fritz. Responsible disclosure of generative models using scalable fingerprinting. In *ICLR*, 2021. 1, 2

- [68] Zhongqi Yue, Tan Wang, Qianru Sun, Xian-Sheng Hua, and Hanwang Zhang. Counterfactual zero-shot and open-set visual recognition. In *CVPR*, pages 15404–15414, 2021. [5](#)
- [69] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deep-fake detection. In *CVPR*, pages 2185–2194, 2021. [1](#)
- [70] Yu Zheng, Yueqi Duan, Jiwen Lu, Jie Zhou, and Qi Tian. Hyperdet3d: Learning a scene-conditioned 3d object detector. In *CVPR*, pages 5585–5594, 2022. [2](#)
- [71] Yu Zheng, Yueqi Duan, Zongtai Li, Jie Zhou, and Jiwen Lu. Learning dynamic scene-conditioned 3d object detectors. *TPAMI*, 46(5):2981–2996, 2023. [2](#)
- [72] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, pages 2921–2929, 2016. [5](#)
- [73] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networkss. In *ICCV*, 2017. [6](#)