

UIAMark: Unified Identity and Attribute Watermarking for Source Tracing and Proactive Deepfake Detection

Peiqi Jiang*, Yuhao Sun*, Lingyun Yu and Hongtao Xie

USTC, EEIS

{jpqjiang, syh3327}@mail.ustc.edu.cn {yuly, htxie}@ustc.edu.cn

Abstract

Malicious Deepfakes have posed sharp security risks to our society by generating remarkably realistic forged faces. Although various countermeasures have been developed for the post-detection of Deepfakes, their performance still lacks reliability and stability in real-world applications. Instead of detecting Deepfakes passively, we propose a novel framework called **Unified Identity and Attribute Watermarks (UIAMark)**, which is devised to provide a proactive strategy for source tracing and Deepfake-agnostic detection with a face disentangling and reconstruction architecture. In detail, we embed a robust watermark into the **facial-irrelevant attributes** which is naturally resilient to Deepfakes, accomplishing robust **source tracing**. Meanwhile, pairs of self-correlated semi-fragile watermarks are injected into **facial-relevant attributes** and **identity** features. If either the facial identity or attributes have been manipulated, the correlation between the semi-fragile watermarks will be broken. In this way, we can fully leverage the inherent Deepfake-fragile property of the high-level semantic representations, achieving **Deepfake-agnostic detection**. Extensive experiments demonstrate the effectiveness of our approach, we achieve an average detection accuracy of 97.36% across various Deepfake methods and a low BER of 0.0389% across different distortions.

1 Introduction

With the rapid evolution of deep generative models [Karras *et al.*, 2019; Goodfellow *et al.*, 2020], advanced face forgery algorithms [Thies *et al.*, 2016; Thies *et al.*, 2019] are able to create remarkably realistic face images, posing potential risks for malicious use and raising significant privacy and security concerns. In recent years, there has been a substantial increase in research dedicated to detecting Deepfakes [Zheng *et al.*, 2021; Haliassos *et al.*, 2021]. Most of these works formulated this task as a binary classification problem and focused on the in-dataset performance. However, training

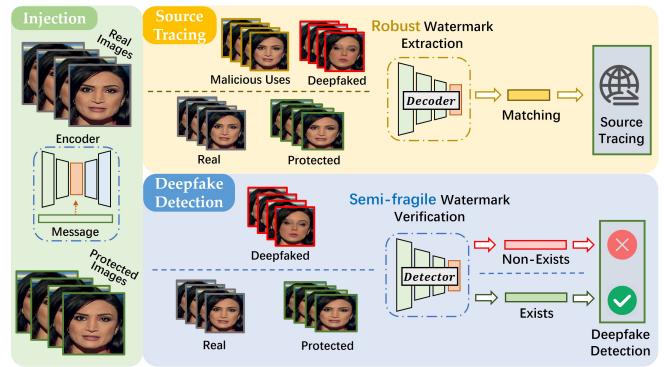


Figure 1: Workflow of source tracing and proactive Deepfake detection. 1) For **Injection**, we embed secret messages into the target image with different robustness levels, and it may be manipulated by malicious users. 2) For **Source Tracing**, we extract the robust watermark and matching with the predefined criteria. 3) For **Deepfake Detection** we verify the existence of the semi-fragile watermark. If the image has been Deepfaked, this watermark will not exist.

a classifier that solely relies on pre-existing forgery artifacts may not be sufficient for detecting unseen manipulations. While several attempts [Shiohara and Yamasaki, 2022; Dong *et al.*, 2023] have been made to improve generalization ability by capturing common forgery clues, nevertheless, their performance is still limited when applied in real-world scenarios.

Alternatively, a distinct area of research has focused on proactive forensics against malicious tampering, which leverages digital watermarking techniques [Zhu *et al.*, 2018] to safeguard real images by embedding secret messages in advance. Based on the robustness of the watermark, we classify these watermarking methods into two distinct types. 1) The *robust watermark* is resilient against both common distortions and Deepfake manipulations, making it suitable for trustworthy source verification and source tracing [Wang *et al.*, 2020]. 2) The *semi-fragile watermark* is robust to common distortions but sensitive to malicious Deepfakes, thus we can proactively detect Deepfakes by checking whether the watermark exists in the tampered image [Asnani *et al.*, 2022].

Although several attempts have been made in proactive forensics, the field continues to grapple with two key challenges: **1) Unrealistic premise of watermark message**

*Equal contribution.

63 **prior:** Most existing methods rely on injecting semi-fragile
64 watermarks, and the detection process involves matching the
65 decoded semi-fragile watermarks with the original water-
66 marks to identify Deepfakes. However, in practical water-
67 mark distribution scenarios, the secret message embedded in
68 different images is not a constant value. Therefore, obtaining
69 prior knowledge of the original watermark during the detec-
70 tion process becomes unrealistic. **2) Lack of Generaliza-**
71 **tion:** Current research fails to achieve Deepfake-agnostic de-
72 tection. For instance, FakeTagger [Wang *et al.*, 2021] uses a
73 GAN simulator to act as a surrogate module to generate Deep-
74 fake translations. However, this strategy can lead to inconsis-
75 tent outcomes when confronted with manipulations beyond
76 the scope of the training process. A similar challenge is evi-
77 dent in [Wu *et al.*, 2023], where their method proves effective
78 against specific types of Deepfake generators but struggles to
79 generalize against others.

80 To tackle the aforementioned issues, we propose a **Unified**
81 **Identity and Attribute Watermarking (UIAMark)** frame-
82 work to provide robust source tracing and self-verifiable
83 Deepfake-agnostic detection. The workflow of our approach
84 is illustrated in Figure 1. To realize the proposed objective,
85 we analyze the watermark embedding process in two critical
86 aspects: 1) For reliable source tracing, we embed a decod-
87 able (i.e., specific binary sequence can be extracted) robust
88 watermark that is resilient against various types of tamper-
89 ing. Moreover, to hide the robust watermark in the facial-
90 irrelevant attributes, we adopt a learning-based encoder-
91 decoder watermarking architecture to embed the watermark
92 adaptively. 2) To achieve proactive detection without the wa-
93 termark message prior, a pair of correlated semi-fragile wa-
94 termarks that can be self-verified are embedded. Moreover, to
95 accomplish Deepfake-agnostic detection, we inject the corre-
96 lated watermark sequences into the facial-relevant attributes
97 and identity features in an additive way. If either the fa-
98 cial identity or attributes are manipulated, it will disrupt the
99 correlation between the semi-fragile watermarks, making our
100 approach effective against both face-swapping and attribute-
101 editing Deepfake manipulations.

102 Building upon these insights, UIAMark embeds two as-
103 pects of watermarks within the facial attributes and identi-
104 ties through a modified face reconstruction network. Specif-
105 ically, given an original image, we first generate a pair of
106 k -bit control sequence and $2^k - 1$ -bit pseudo-random noise
107 (PN) sequence with unique autocorrelation properties tied to
108 each specific control sequence. Then we embed the binary
109 control sequence into the face attributes in a deep separable
110 manner, where a single encoder is used for injecting and two
111 decoders are used to extract the secret message at different
112 robustness levels. In this way, we can obtain the robust wa-
113 termark on facial-irrelevant attributes for **source tracing** and
114 a semi-fragile watermark on facial-relevant attributes which
115 are inherently fragile to attribute-edit manipulations. Mean-
116 while, the identity watermark is injected by entangling the
117 facial identity feature with the PN sequence, making the iden-
118 tity watermark inherently fragile to face-swap translations.
119 Therefore, we can accomplish the **Deepfake-agnostic de-**
120 **tection** by verifying the correlation status between the semi-
121 fragile attribute watermark and identity watermark, any tam-

122 pering with either facial attributes or identity will disrupt the
123 inherent auto-correlation property of the PN sequence. Fi-
124 nally, based on the watermarked attributes and watermarked
125 identity, we leverage a face generator to generate the final
126 watermarked image. To conclude, our contributions can be
127 summarized as:

- We propose a unified attribute and identity watermark-
128 ing framework, which brings a new paradigm to source
129 tracing and proactive Deepfake detection.
- We accomplish Deepfake-agnostic detection by leverag-
130 ing the inherent fragility of high-level facial semantic
131 information (i.e., facial-relevant attributes and identity
132 features) and utilizing the correlation property between
133 the control and pseudo-random sequences.
- Extensive experiments demonstrate the effectiveness of
134 our approach in source tracing and proactive Deepfake
135 detection. In particular, we achieve a BER of 0.0389%
136 across different distortions and an average detection ac-
137 curacy of 97.36% across various Deepfake methods in
138 the black-box setting.

2 Related Work

2.1 Passive Deepfake Detection

With the rise of malicious Deepfake abuses, numerous at-
144 tempts have been in the field of passive Deepfake detection.
145 These methods all follow a similar paradigm that is trained on
146 a given dataset with binary labels and detect Deepfakes after
147 the manipulations have already happened. Earlier research
148 primarily focused on designing a powerful network to detect
149 artifacts from RGB representations [Amerini *et al.*, 2019].
150 Besides, researchers have also begun to notice the vital prob-
151 lem of cross-dataset generalization. For instance, Face X-
152 ray [Li *et al.*, 2020] focused on detecting blending boundary
153 artifacts, PCL+I2G [Zhao *et al.*, 2021] detecting Deepfakes
154 by capturing the intra-frame inconsistency, and SBIs [Shio-
155 hara and Yamasaki, 2022] reproduce common Deepfake arti-
156 facts by seamlessly stitching two images.

2.2 Proactive Defence against Deepfake

The research on proactive defense mainly leverages digital
159 watermarking techniques to safeguard real images by em-
160 bedding secret messages in advance. [Wang *et al.*, 2021]
161 proposed a Deepfake-robust watermarking network to iden-
162 tify the Deepfake source. [Asnani *et al.*, 2022] introduced
163 a semi-fragile watermarking strategy through a combination
164 of Deepfake surrogate models. [Zhao *et al.*, 2023] lever-
165 aged identity watermarking to detect face-swap manipula-
166 tions. [Wu *et al.*, 2023] proposed a separable deep water-
167 marking approach that contains both robust watermark and
168 semi-fragile watermark, and widens the application scenarios
169 of watermark protection methods. However, most prevalent
170 methods rely on conventional digital watermarking frame-
171 works and substitute the noise layer with a Deepfake surro-
172 gate module in their design. This training approach inher-
173 ently restricts the capacity for achieving generalizable Deep-
174 fake detection. In contrast, our framework is devised from a

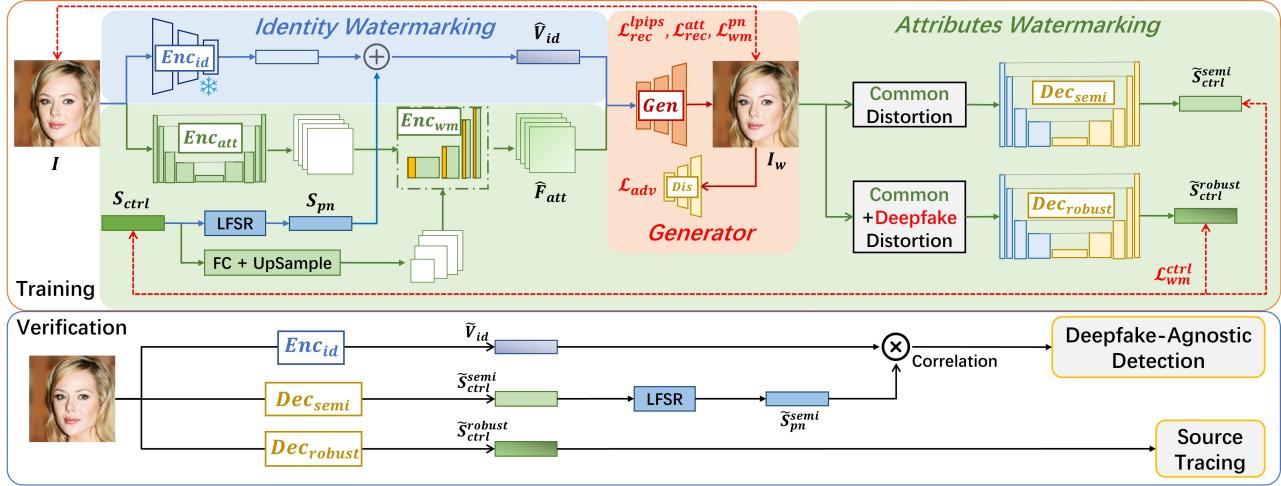


Figure 2: The overview of UIAMark. 1) In the training stage, we disentangle the input image into the attributes and identity representations. Training. Then, we embed the control sequence into attribute features through a watermark encoder, and it can be separately extracted by two decoders at different robustness levels. Meanwhile, we inject the PN sequence into the identity feature, making the identity watermark inherently fragile to face-swap translations. 2) In the verification stage, we leverage the robust watermark for **source tracing** and **detect Deepfakes** based on the correlation result between the identity features and the control sequence decoded from the semi-fragile watermark.

176 face reconstruction network, allowing us to leverage the inherent fragility of high-level facial semantic information for
177 Deepfake-agnostic detection.
178

179 2.3 Digital Watermarking

180 Digital watermarking is a broadly studied approach to im-
181 perceptibly concealing secret messages within media con-
182 tent [Cox, 2007]. However, earlier hand-crafted watermark-
183 ing methods suffer from limited robustness and generalization.
184 Recently, DNN-based approaches have been introduced
185 in digital watermarking for stronger robustness. For example,
186 HiDDeN [Zhu *et al.*, 2018] is the first end-to-end framework
187 for robust watermarking, where an encoder-decoder network
188 is utilized for watermark injection and extraction, along with
189 a noise layer to simulate perturbations for robustness. [Luo
190 *et al.*, 2020] employs an adversarial noise layer to further
191 improve the robustness against unseen distortions. Inspired
192 by the encoder-decoder-based watermarking paradigm, our
193 method designs a single encoder and two separate decoders
194 which are integrated into a face reconstruction network for
195 the attribute watermarking.

196 3 Method

197 3.1 Model Architecture

198 We illustrate the overview of the proposed UIAMark in Figure 199 2. The whole pipeline is based on a face reconstruction
200 network and we elaborate on the watermark injection and
201 extraction process as follows: 1) Given an real image I , we
202 first randomly generate a pair of k -bit control sequence S_{ctrl}
203 and $2^k - 1$ -bit PN sequence S_{pn} as the secret messages. 2)
204 Then, we use an identity encoder and an attribute decoder to
205 disentangle the input image into two independent represen-
206 tations: attributes and identity. 3) Next, S_{ctrl} and S_{pn} are
207 embedded into the attributes and identity separately. For the

208 attribute watermark, we leverage a learnable watermark en-
209 coder to embed the control sequence S_{ctrl} into the attributes
210 and obtain the watermarked attributes. For the identity water-
211 mark, we directly add the scaled S_{pn} into the identity vector
212 to generate the watermarked identity. 4) Finally, we utilize a
213 face reconstruction generator to synthesize the correspond-
214 ing watermarked image I_w .

215 In the watermark extraction stage, we first utilize two sep-
216 arate watermark decoders to extract the attribute watermark
217 at two robustness levels, i.e., a robust attribute watermark
218 S_{ctrl}^{robust} and a semi-fragile attribute watermark S_{ctrl}^{semi} , where
219 the robust watermark can be used for **source-tracing**. Mean-
220 while, we extract the identity vector from the watermarked
221 image and compute the correlation with the S_{ctrl}^{semi} sequence,
222 if Deepfakes manipulate either attribute or identity, the corre-
223 lation property will break. Thus, we can leverage this feature
224 to **detect Deepfakes**.

225 Next, we will introduce the details of each component.
226 **Sequence Generation.** In the sequence generation stage, we
227 aim to create a pair of messages with certain correlations
228 for the attribute and identity watermarking, therefore we can
229 leverage this correlation to verify the existence of the iden-
230 tity watermark. Inspired by the telecommunication system,
231 we fully leverage the features of pseudo-random noise (PN)
232 sequences. We first generate a random k -bit binary sequence
233 denoted as the control sequence S_{ctrl} . Next, we feed the con-
234 trol sequence into the linear feedback shift registers (LFSR)
235 to generate a $2^k - 1$ -bit PN sequence S_{pn} . The PN sequence
236 has the same statistical property as the Gaussian noise but
237 has a special correlation property. The auto-correlation of the
238 PN sequence is approximate to an impulse function, with the
239 impulse peaks occurring only at the zero point, it can be for-
240 mulated as:

$$Corr[l] = \begin{cases} 1 & \text{if } l = 0 \\ -\frac{1}{N} & \text{otherwise} \end{cases}, \quad (1)$$

where the period N equals to $2^k - 1$. With this autocorrelation property, we can verify whether the control sequence and the PN sequence match with each other.

Feature Disentanglement. To entangle the message with the original facial attributes and identity, we first employ two encoders, namely identity encoder $\text{Enc}_{\text{id}}(\cdot)$ and attributes encoder $\text{Enc}_{\text{att}}(\cdot)$, to extract the respective representations from the input image I .

The attributes representation of a face image encompasses features including style, pose, expression, background, etc. These features contain different levels of detail, ranging from coarse features like the overall facial outline to finer details such as hair color. Therefore, we devise a modified U-net style network to encode k levels of attribute features $F_{\text{att}} = \{F_{\text{att}}^1, F_{\text{att}}^2, \dots, F_{\text{att}}^n\} = \text{Enc}_{\text{att}}(I)$, where $F_{\text{att}}^i \in \mathbb{R}^{C_i \times H_i \times W_i}$ denotes the i -th level of feature map in the encoder $\text{Enc}_{\text{att}}(\cdot)$.

The identity representation refers to the high-level semantic feature used to characterize a specific person. We utilize a pretrained face recognition model to act as the identity encoder. We extract the identity $V_{\text{id}} = \text{Enc}_{\text{id}}(I)$, where the identity vector $V_{\text{id}} \in \mathbb{R}^l$, and l is the vector length.

During the training process, the attribute encoder $\text{Enc}_{\text{att}}(\cdot)$ is trained from scratch and we freeze the pretrained identity encoder $\text{Enc}_{\text{id}}(\cdot)$. Thus, the disentanglement process does not require extra attribute annotations.

Watermark Encoding. Given the disentangled facial representations, we aim to embed a pair of control sequence S_{ctrl} and PN sequence S_{pn} into the attribute features F_{att} and identity vector V_{id} .

Attributes Watermarking: Considering that attribute features have different detail levels, and their feature maps are at different sizes, we adopt an encoder-decoder-based watermark embedding architecture and aim to inject two robustness levels of watermarks on different attribute spaces. Specifically, for each level of attribute features, the attribute watermark encoder Enc_{wm} maps the k -bit control sequence S_{ctrl} into the corresponding feature map size through independent $FC + \text{Upsample}$ layers. Then we concatenate the upsampled message with F_{att}^i and fuse them through a light-weighted Coordinate Attention [Hou *et al.*, 2021] mechanism, which makes the message embedding process become position-sensitive and direction-aware. The process is illustrated as $\widehat{F}_{\text{att}} = \{\widehat{F}_{\text{att}}^1, \widehat{F}_{\text{att}}^2, \dots, \widehat{F}_{\text{att}}^n\} = \text{Enc}_{\text{wm}}(F_{\text{att}})$, where \widehat{F}_{att} is the watermarked attributes with n levels of watermarked feature maps $\widehat{F}_{\text{att}}^i$ at different resolution.

Identity Watermarking: Observing that facial identity vectors exhibit statistical characteristics similar to Gaussian noise, we inject autocorrelated pseudo-noise (PN) sequence S_{pn} with the same statistical properties into the identity vector, as formalized as $\widehat{V}_{\text{id}} = V_{\text{id}} + \gamma S_{\text{pn}}$, where \widehat{V}_{id} is the watermarked identity and γ is the scale factor.

Face Reconstruction. To integrate watermarked attributes and identity and generate watermarked images with high visual quality. We employ the widely used *Adaptively Attentional Denormalization Generator* (AAD-Generator) [Li *et al.*, 2019] in the face-swapping domain for facial reconstruc-

tion. Specifically, the AAD-Generator contains a cascade of AAD Residual Blocks, and each block consists of two AAD layers with residue convolution connections. AAD layer first takes the normalized activation from the previous level as input h_{in}^i . The attentional mask M_i is generated from h_i with convolutions and a Sigmoid operation. The attribute and identity are integrated by denormalizing h_i , it can be formulated as $A_i = \gamma_{\text{att}}^i \otimes h_{\text{in}}^i + \beta_{\text{att}}^i$ and $I_i = \gamma_{\text{id}}^i \otimes h_{\text{in}}^i + \beta_{\text{id}}^i$, where A_i and I_i is the integrated attribute and identity in the i -th layer, γ_{att}^i and β_{att}^i are two modulation parameters both convolved from \widehat{F}_{att} , they share the same tensor dimensions with $h_{\text{in}}^i \in \mathbb{R}^{C_i \times H_i \times W_i}$. Similarly, $\gamma_{\text{id}}^i \in \mathbb{R}^{C_i}$ and $\beta_{\text{id}}^i \in \mathbb{R}^{C_i}$ are another two modulation parameters generated from \widehat{V}_{id} through FC layers.

Finally, after the integration of n level of attribute and identity features, the AAD-Generator $\text{Gen}_{\text{AAD}}(\cdot)$ outputs the final output image as $I_w = \text{Gen}_{\text{AAD}}(\widehat{F}_{\text{att}}, \widehat{V}_{\text{id}})$.

Watermark Extraction. Given a watermarked image I_w , the decoding processes for facial attributes and identity watermarks differ based on their respective embedding methods.

Attribute Watermark Decoding: For the attribute watermark, considering that different levels of attribute features include coarse-grained background and contours, as well as finer details such as hair color and facial expressions, we opted for a separable watermark design [Wu *et al.*, 2023] to decode two watermarks at different robustness levels. We extract a robust watermark and a semi-fragile watermark. The decoding process is formulated as

$$\begin{cases} \tilde{S}_{\text{ctrl}}^{\text{robust}} = \text{Dec}_{\text{robust}}(\text{Noise}(\text{DF}(I_w))) \\ \tilde{S}_{\text{ctrl}}^{\text{semi}} = \text{Dec}_{\text{semi}}(\text{Noise}(I_w)) \end{cases}, \quad (2)$$

where $\tilde{S}_{\text{ctrl}}^{\text{robust}}$ and $\tilde{S}_{\text{ctrl}}^{\text{semi}}$ denotes the decoded robust and semi-fragile attribute watermarks. $\text{Noise}(\cdot)$ denotes a differential noise layer containing JPEG compression, resize, crop, and random noise distortions. $\text{DF}(\cdot)$ denotes surrogate Deepfake models for the robust watermark, including SimSwap [Chen *et al.*, 2020] and StarGAN [Choi *et al.*, 2018].

Identity Watermark Extracting: During decoding, it is necessary to detect whether the corresponding sequence exists in the identity features of the watermarked image (further details on the verification method are presented in Section 3.3). We extracted the reconstructed identity vector as $\widetilde{V}_{\text{id}} = \text{Enc}_{\text{id}}(I_w)$, where $\widetilde{V}_{\text{id}}$ denotes the identity of the watermarked image.

3.2 Training Loss

Message Loss. For the robust and semi-fragile attribute watermarks, we use the L_2 loss to constrain the message before and after decoding. The formula is as follows:

$$\mathcal{L}_{\text{wm}}^{\text{ctrl}} = L_2(S_{\text{ctrl}}, \tilde{S}_{\text{ctrl}}^{\text{robust}}) + L_2(S_{\text{ctrl}}, \tilde{S}_{\text{ctrl}}^{\text{semi}}). \quad (3)$$

For the identity watermark, we do not decode specific watermark values. Instead, we constrain the cosine similarity between the identity vector in the watermarked image and the original identity watermark vector. The process is as follows:

$$\mathcal{L}_{\text{wm}}^{\text{pn}} = 1 - \cos(\widetilde{V}_{\text{id}} - \widehat{V}_{\text{id}}). \quad (4)$$

347 **Reconstruction Loss.** To maintain the high fidelity of the re-
 348 reconstructed image to the original and minimize conflicts with
 349 watermark injection at the pixel level, we utilize a perceptual
 350 similarity loss (LPIPS) [Zhang *et al.*, 2018] instead of
 351 the pixel-level reconstruction loss:

$$\mathcal{L}_{rec}^{lpipl} = \|\text{LPIPS}(I, I_w)\|_2. \quad (5)$$

352 Additionally, to ensure the training stability of the face recon-
 353 struction model, we also constrain the reconstruction loss for
 354 attributes as follows:

$$\mathcal{L}_{rec}^{att} = \frac{1}{2} \sum_{i=1}^n \left\| \widehat{F}_{att}^i - \widetilde{F}_{att}^i \right\|_2^2, \quad (6)$$

355 where \widetilde{F}_{att}^i is the attribute features extracted from the water-
 356 marked image I_w .

357 **Adversary Loss.** To enhance the realism of the recon-
 358 structed image, we utilize a multi-scale discriminator $\text{Dis}(\cdot)$
 359 from [Isola *et al.*, 2017] with hinge loss functions, training
 360 our model in an adversarial manner:

$$\mathcal{L}_{adv} = \log(\text{Dis}(I)) + \log(1 - \text{Dis}(I_w)). \quad (7)$$

361 **Overall Loss.** Our framework is trained using a weighted
 362 sum of the aforementioned losses, as defined by:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{wm}^{ctrl} + \lambda_2 \mathcal{L}_{wm}^{pn} + \lambda_3 \mathcal{L}_{rec}^{lpipl} + \lambda_4 \mathcal{L}_{rec}^{att} + \lambda_5 \mathcal{L}_{adv}, \quad (8)$$

363 where $\lambda_1 \sim \lambda_5$ are hyper-parameters to weight the corre-
 364 sponding loss.

365 3.3 Watermark Verification

366 In this section, we will elaborate on how to use attribute and
 367 identity watermarks for source tracing and Deepfake detec-
 368 tion. Source tracing requires the secret message from a robust
 369 watermark and cross-verification with the original message.
 370 Thus, we leverage the robust attribute watermark $\widetilde{S}_{ctrl}^{robust}$ to
 371 accomplish it. As for Deepfake detection, we aim to detect
 372 the presence of the PN sequence S_{pn} in the identity fea-
 373 ture of the watermarked image. The existence of S_{pn} will
 374 determine whether the image is Deepfaked or not. Given
 375 the watermarked image is not modified by Deepfakes, we
 376 extract its identity vector \widetilde{V}_{id} , and it can be expanded as
 377 $\widetilde{V}_{id} = V_{id} + \gamma S_{pn}$, where V_{id} is the clean identity vector,
 378 S_{pn} is the PN sequence. Then, we compute its correlation
 379 with the PN sequence decoded from the semi-fragile water-
 380 mark $\widetilde{S}_{ctrl}^{semi}$ to generate a corresponding PN sequence $\widetilde{S}_{pn}^{semi}$
 381 through the LFSR. Then verify whether the watermark exists
 382 in the input image. The correlation process is defined as:

$$\begin{aligned} Corr[l] &= \sum_{i=0}^{N-1} \widetilde{S}_{pn}^{semi}[i] * \widetilde{V}_{id}[j] \\ &= \sum_{i=0}^{N-1} \left(\widetilde{S}_{pn}^{semi}[i] * \gamma S_{pn}[j] + \widetilde{S}_{pn}^{semi}[i] * V_{id}[j] \right), \end{aligned} \quad (9)$$

383 where $j = (N - 1 + i - l)$, N is the sequence length and
 384 l is ranged from $\{-N + 1, \dots, 0, \dots, N - 1\}$. The first term
 385 in Equation (9) is equivalent to an impulse function and the

second term is a bounded constant. Thus, according to Equa-
 386 tion (1), $Corr[l]$ will have a distinct peak value at 0 index if
 387 the watermarked image is not manipulated. Thus we can de-
 388 termine whether the watermarked image has been Deepfaked
 389 by comparing the Peak to Average Power Ratio (PAPR) with
 390 threshold τ .
 391

4 Experiments

392 4.1 Implementation Details

393 **Experiments Settings.** Our experiments are conducted on
 394 the CelebA-HQ dataset [Karras *et al.*, 2017] following the
 395 official split for training, validation, and testing. We further
 396 conduct experiments on the Flickr-Faces-HQ (FFHQ) [Karras
 397 *et al.*, 2019] dataset to provide more comprehensive results.
 398 Unless stated otherwise, all images in the experiment have
 399 been aligned and cropped to the size of 256×256 . We set
 400 $\lambda_1 = \lambda_2 = 1$, $\lambda_3 = \lambda_4 = 10$, $\lambda_5 = 0.1$, scale factor
 401 $\gamma = 0.1$, and threshold $\tau = 4$. As for the choice of PN
 402 sequence type, we select the widely used M Sequence.
 403

404 **Evaluation Metrics.** For Deepfake detection, we compute
 405 image-level Accuracy (Acc) and F1-Score. For source trac-
 406 ing, we evaluate the robust watermarks with the average bit
 407 error rate (BER) metric. For visual quality, we report the
 408 average PSNR, SSIM, and LPIPS for watermarked images
 409 through the testing set.
 410

411 **Compared Methods.** We mainly compared our work with
 412 SepMark [Wu *et al.*, 2023], which is the first work of unified
 413 source tracing and Deepfake detection. Additionally, most
 414 existing works only embed watermarks with a single func-
 415 tionality. Therefore, we adopt robust watermarking meth-
 416 ods like HiDDeN [Zhu *et al.*, 2018], MBRS [Jia *et al.*,
 417 2021], and semi-fragile watermarking methods such as Face-
 418 Sign [Neekhara *et al.*, 2022], IDMark [Zhao *et al.*, 2023] as
 419 our baselines. Note that SepMark only reports the BER of ro-
 420 bust and semi-fragile watermarks, thus we compare the delta
 421 of two watermarks' BER with a certain threshold to deter-
 422 mine whether the image is Deepfaked (See details in supple-
 423 mentary materials). We also compare the proposed proac-
 424 tive defense strategy with SOTA detectors such as ICT [Dong
 425 *et al.*, 2022], SBI [Shiohara and Yamasaki, 2022], and UIA-
 426 ViT [Zhuang *et al.*, 2022].
 427

4.2 Visual Quality

428 We evaluate the visual quality of the generated watermark im-
 429 ages, as shown in Table 1. Our method achieves competitive
 430 results against the SepMark method (33.43 of PNSR, 0.9446
 431 of SSIM, and 0.0021 of LPIPS). Additionally, as illustrated in
 432 Figure 3, no artifacts can be observed by the naked eye in the
 433 generated watermark images, demonstrating the high fidelity
 434 of our approach. Specifically, from the residual between the
 435 watermark image and the original image, it can be observed
 436 that in regions associated with facial identity, the distribution
 437 of watermarks is similar to white noises. Conversely, in re-
 438 gions related to the background and certain facial attributes,
 439 the watermark is distributed differently in three RGB chan-
 440 nels. This demonstrates the effectiveness of our approach
 441 to separately embedding watermarks on attribute and identity
 442 information.
 443

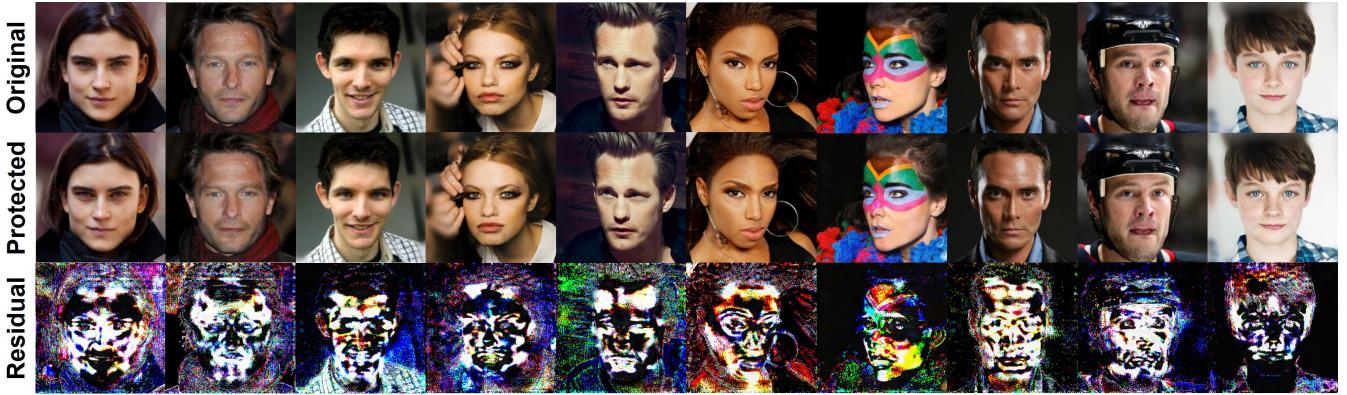


Figure 3: Visual Quality of watermarked images. From top to bottom are the original image, watermarked image, and normalized residual between the original and watermarked images.

Method	CelebA-HQ			FFHQ		
	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS
HiDDeN	33.26	0.8880	0.0210	31.13	0.8670	0.0170
FaceSign	32.33	0.9211	0.0260	32.99	0.8890	0.0150
MBRS	33.05	0.8106	0.0141	33.00	0.7750	0.0241
IDMark	33.32	0.9400	0.0030	33.50	0.9500	0.0042
SepMark	38.56	0.9328	0.0080	37.60	0.8793	0.0569
Ours	<u>33.43</u>	0.9446	0.0021	<u>33.75</u>	0.9533	0.0023

Table 1: Visual quality of the watermarked image. We achieve the best fidelity on SSIM and LPIPS and the second-best on PSNR.

4.3 Source Tracing

The performance of sourcing tracing primarily depends on the robust attribute watermark $\tilde{S}_{ctrl}^{robust}$, which should have a low BER under both common and Deepfake distortions. As shown in Table 4, the average BER of $\tilde{S}_{ctrl}^{robust}$ in our method achieves 0.0389% under common distortions and 0.1766% under Deepfake distortions, demonstrating stronger robustness compared to the baselines. This verifies the advantage of our hierarchical fine-grained attribute watermark embedding and separable design. The $\tilde{S}_{ctrl}^{robust}$, which is mainly distributed in the background region, can better handle strong distortions, especially background-independent Deepfake distortions. Notably, due to the learning-based attribute watermarking paradigm, a portion of $\tilde{S}_{ctrl}^{robust}$ is inevitably dispersed across facial-relevant attributes. Consequently, when encountering face-swapping techniques such as SimSwap, there may be a marginal increase in BER.

4.4 Deepfake Detection

For Deepfake detection, both SOTA detector-based and watermark-based methods are selected for comparison with our method. We evaluate these methods on four different Deepfake models including face-swap methods like SimSwap and DiffFace, as well as attribute-edit methods like StarGAN and HFGI. As shown in Table 3, our method achieves satisfactory Deepfake detection accuracy across all the models with an average of 97.36%. It is worth mentioning that Sep-

Distortion	Robust Watermark BER(%) ↓			
	FaceSign	MBRS	SepMark	Ours
Identity	0.0136	0.0000	0.0000	0.0000
JPEG	0.8258	0.2597	0.2136	0.0782
Resize	1.0726	0.0000	0.0059	0.0000
GaussianBlur	0.1671	0.0000	0.0024	0.0010
MedianBlur	0.0977	0.0000	0.0012	0.0100
Brightness	10.8196	0.0000	0.0059	0.0079
Contrast	0.0334	0.0000	0.0012	0.0011
Saturation	0.7116	0.0000	0.0000	0.0000
Hue	8.3780	0.0000	0.0000	0.0000
SaltPepper	12.3238	0.0000	0.0413	0.0177
GaussianNoise	7.0697	0.0000	0.7460	0.3120
Avg. Common	3.7738	0.0236	0.0848	<u>0.0389</u>
SimSwap	0.0537	19.3744	13.8255	0.2231
StarGAN	0.5311	18.1358	0.1268	0.1300
Avg. Malicious	<u>0.2924</u>	18.7551	6.9762	0.1766

Table 2: Quantitative comparison between SepMark, IDMark, and ours on CelebA-HQ regarding the bit-wise error rate under common and malicious Deepfake manipulations.

Mark performs slightly better than our method in StarGAN. This is because SepMark incorporates StarGAN and SimSwap during the training process, whereas for our method, they are considered as black-box scenarios. Conversely, on the unseen Deepfake methods like HFGI and DiffFace, the detection accuracy of SepMark significantly drops to nearly 50%, while our method maintains high accuracy. The results demonstrate the superiority of our method in Deepfake-agnostic detection by leveraging the inherent fragility of both the identify and attribute information.

Meanwhile, we evaluate the performance of our method in Deepfake detection under common distortions, as shown in Table 4. Our method maintains a high detection accuracy under different distortions, verifying the robustness of the semi-fragile watermarks embedded in identity and attribute information against common distortions.

Method	CelebA-HQ (Acc(%) ↑ / F1-Score(%) ↑)				FFHQ (Acc(%) ↑ / F1-Score(%) ↑)			
	SimSwap	StarGAN	HFGI	DiffFace	SimSwap	StarGAN	HFGI	DiffFace
ICT	76.70/76.63	77.90/78.56	69.72/68.19	76.80/78.08	76.95/76.49	78.80/78.82	68.15/67.74	78.45/78.51
SBI	76.65/77.08	70.90/76.34	70.48/75.91	42.96/32.34	71.83/76.46	74.53/79.18	71.53/76.15	44.86/37.49
UIA-ViT	48.65/19.07	62.21/50.94	71.33/66.71	48.95/19.92	48.75/65.25	69.21/66.36	69.01/66.07	45.90/26.07
IDMark	<u>94.05/94.02</u>	78.95/81.72	<u>89.35/89.87</u>	<u>95.00/94.92</u>	92.80/92.58	80.00/81.63	<u>90.70/88.07</u>	<u>91.83/91.57</u>
SepMark	89.67/90.63 [†]	99.78/99.78[†]	50.00/66.67	50.00/66.67	87.60/88.71 [†]	98.73/98.73[†]	50.00/66.67	50.00/66.67
Ours	97.42/97.48	<u>97.50/96.97</u>	97.28/97.23	97.25/97.19	98.90/98.89	<u>96.90/96.94</u>	98.67/98.66	98.07/98.07

Table 3: Performance of Deepfake detection Accuracy ↑ and F1-score ↑. We compare the proposed proactive defense strategy with SOTA detectors, as well as watermarking baselines. †: White-box evaluation.

Distortion	Acc(%) ↑			
	SimSwap	StarGAN	HFGI	DiffFace
Identity	97.42	97.50	97.28	97.25
JPEG	93.63	76.75	87.95	77.42
Resize	97.15	97.00	97.20	97.58
GaussianBlur	90.52	89.95	89.80	97.95
MedianBlur	93.23	92.80	93.75	97.65
Brightness	89.97	87.45	76.45	86.67
Contrast	94.52	94.05	94.95	96.27
Saturation	97.13	96.70	97.45	97.93
Hue	96.85	96.85	97.55	97.20
SaltPepper	83.53	72.75	74.05	71.37
GaussianNoise	92.37	77.90	79.35	73.83
Avg.	93.30	89.06	89.62	90.10

Table 4: Robustness of Deepfake detection on CelebA-HQ dataset. Our method maintains a high accuracy under various distortions.

γ	Visual Quality		Acc(%) ↑			
	PSNR	SSIM	SimSwap	StarGAN	HFGI	DiffFace
0.01	34.75	0.9501	53.17	52.18	52.83	52.67
0.05	34.11	0.9480	85.67	86.15	83.00	84.73
0.10	33.43	0.9446	97.42	97.50	97.28	97.25
0.25	32.08	0.9400	97.67	97.83	94.67	94.83
0.50	30.63	0.9340	89.83	89.33	83.17	86.63
1.00	26.87	0.9150	53.17	50.83	51.83	50.13

Table 5: Ablations of scale factor γ in identity watermark. We find that $\gamma = 0.1$ is the best trade-off point.

4.5 Ablation Study

Impact of scale factor γ in identity watermark. We investigate the impact of different scale factors on identity watermarks. The ideal injection strength needs to be strong enough for detection while maintaining high fidelity. As illustrated in Table 5, we find that $\gamma = 0.1$ for the best trade-off point.

Impact of different identity watermark sequences. According to Equation (9), the correlation property of the injected identity watermark sequence plays a significant role in the verification. To analyze the impact of the different sequence types in the identity watermarking process, we select two representative PN sequences: M Sequence and Gold Code. Along with two common random sequences: Bernoulli and Gaussian Sequences to compare the detection perfor-

ID Watermark Sequence Type	Acc(%) ↑			
	SimSwap	StarGAN	HFGI	DiffFace
Random Bernoulli	50.00	52.10	53.21	51.12
Random Guassain	51.20	53.25	52.42	52.50
Pseudo Glod Code	<u>96.56</u>	97.72	<u>96.31</u>	<u>97.00</u>
Random M Sequence	97.42	<u>97.50</u>	97.28	97.25

Table 6: Different sequences' detection results. Demonstrating the necessity of choosing the PN sequences for identity watermarking.

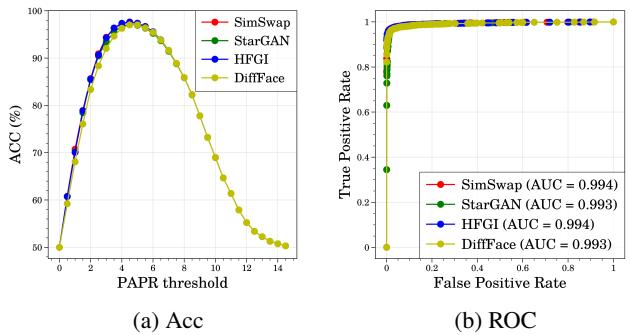


Figure 4: Ablations of PAPR threshold τ , we compute the Acc-PAPR curve and the ROC curve.

mance. As illustrated in Table 6, two PN sequences achieve competitive detection accuracy, while the random sequences failed to detect Deepfakes, demonstrating the necessity of choosing the PN sequences for identity watermarking.

Impact of different PAPR threshold τ . Figure 4a illustrates the impact of different PAPR thresholds τ on the accuracy of Deepfake detection. It can be observed that the optimal threshold for various Deepfake models is approximately equal to 4, which aligns with the value chosen in our experiments. With the variation of threshold values, we also plot the ROC curves and compute corresponding AUC as shown in Figure 4b. For all four Deepfake models, the AUC values are as high as 0.994. These results indicate the excellent performance of our method in Deepfake detection.

5 Conclusion

In this paper, we proposed a novel watermarking paradigm for protecting face images from malicious Deepfakes. By em-

498
499
500
501
502
503
504
505
506
507
508
509
510
511

512
513
514

515 bedding pairs of self-correlated watermarks into the facial at-
516 tributes and identity, we fully leverage the inherent fragility of
517 high-level facial semantic information, accomplishing unified
518 robust source tracing and Deepfake-agnostic detection. Ex-
519 tensive experiments demonstrate our method’s superior per-
520 formance on widely used datasets, establishing a strong base-
521 line for future research.

References

- [Amerini *et al.*, 2019] Irene Amerini, Leonardo Galteri, Roberto Caldelli, and Alberto Del Bimbo. Deepfake video detection through optical flow based cnn. In *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pages 0–0, 2019. 522
- [Asnani *et al.*, 2022] Vishal Asnani, Xi Yin, Tal Hassner, Si-
jia Liu, and Xiaoming Liu. Proactive image manipulation
detection. In *Proceedings of the IEEE/CVF Conference on
Computer Vision and Pattern Recognition*, pages 15386–
15395, 2022. 528
- [Chen *et al.*, 2020] Renwang Chen, Xuanhong Chen, Bing-
bing Ni, and Yanhao Ge. Simswap: An efficient frame-
work for high fidelity face swapping. In *Proceedings of the
28th ACM International Conference on Multimedia*, pages
2003–2011, 2020. 533
- [Choi *et al.*, 2018] Yunjey Choi, Minje Choi, Munyoung
Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Star-
gan: Unified generative adversarial networks for multi-
domain image-to-image translation. In *Proceedings of the
IEEE conference on computer vision and pattern recogni-
tion*, pages 8789–8797, 2018. 538
- [Cox, 2007] IJ Cox. Digital watermarking and steganogra-
phy. *Morgan Kaufmann google schola*, 2:893–914, 2007. 544
- [Dong *et al.*, 2022] Xiaoyi Dong, Jianmin Bao, Dongdong
Chen, Ting Zhang, Weiming Zhang, Nenghai Yu, Dong
Chen, Fang Wen, and Baining Guo. Protecting celebri-
ties from deepfake with identity consistency transformer,
2022. 546
- [Dong *et al.*, 2023] Shichao Dong, Jin Wang, Renhe Ji, Jia-
jun Liang, Haoqiang Fan, and Zheng Ge. Implicit identity
leakage: The stumbling block to improving deepfake de-
tection generalization, 2023. 551
- [Goodfellow *et al.*, 2020] Ian Goodfellow, Jean Pouget-
Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley,
Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Gener-
ative adversarial networks. *Communications of the
ACM*, 63(11):139–144, 2020. 555
- [Haliassos *et al.*, 2021] Alexandros Haliassos, Konstantinos
Vougioukas, Stavros Petridis, and Maja Pantic. Lips don’t
lie: A generalisable and robust approach to face forgery
detection. In *Proceedings of the IEEE/CVF conference on
computer vision and pattern recognition*, pages 5039–
5049, 2021. 560
- [Hou *et al.*, 2021] Qibin Hou, Daquan Zhou, and Jiashi
Feng. Coordinate attention for efficient mobile network
design. In *Proceedings of the IEEE/CVF conference on
computer vision and pattern recognition*, pages 13713–
13722, 2021. 566
- [Isola *et al.*, 2017] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou,
and Alexei A Efros. Image-to-image translation with con-
ditional adversarial networks. In *Proceedings of the IEEE
conference on computer vision and pattern recognition*,
pages 1125–1134, 2017. 571

- 576 [Jia *et al.*, 2021] Zhaoyang Jia, Han Fang, and Weiming
 577 Zhang. Mbrs: Enhancing robustness of dnn-based water-
 578 marking by mini-batch of real and simulated jpeg com-
 579 pression. In *Proceedings of the 29th ACM international*
 580 *conference on multimedia*, pages 41–49, 2021.
- 581 [Karras *et al.*, 2017] Tero Karras, Timo Aila, Samuli Laine,
 582 and Jaakko Lehtinen. Progressive growing of gans for
 583 improved quality, stability, and variation. *arXiv preprint*
 584 *arXiv:1710.10196*, 2017.
- 585 [Karras *et al.*, 2019] Tero Karras, Samuli Laine, and Timo
 586 Aila. A style-based generator architecture for generative
 587 adversarial networks. In *Proceedings of the IEEE/CVF*
 588 *conference on computer vision and pattern recognition*,
 589 pages 4401–4410, 2019.
- 590 [Li *et al.*, 2019] L Li, J Bao, H Yang, D Chen, and
 591 F Wen. Faceshifter: Towards high fidelity and occlu-
 592 sion aware face swapping. arxiv 2019. *arXiv preprint*
 593 *arXiv:1912.13457*, 2019.
- 594 [Li *et al.*, 2020] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao
 595 Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-
 596 ray for more general face forgery detection. In *Proceed-
 597 ings of the IEEE/CVF conference on computer vision and*
 598 *pattern recognition*, pages 5001–5010, 2020.
- 599 [Luo *et al.*, 2020] Xiyang Luo, Ruohan Zhan, Huiwen
 600 Chang, Feng Yang, and Peyman Milanfar. Distortion
 601 agnostic deep watermarking. In *Proceedings of the*
 602 *IEEE/CVF conference on computer vision and pattern*
 603 *recognition*, pages 13548–13557, 2020.
- 604 [Neekhara *et al.*, 2022] Paarth Neekhara, Shehzeen Hussain,
 605 Xinqiao Zhang, Ke Huang, Julian McAuley, and Farinaz
 606 Koushanfar. Facesigns: semi-fragile neural watermarks
 607 for media authentication and countering deepfakes. *arXiv*
 608 *preprint arXiv:2204.01960*, 2022.
- 609 [Shiohara and Yamasaki, 2022] Kaede Shiohara and Toshi-
 610 hiko Yamasaki. Detecting deepfakes with self-blended
 611 images. In *Proceedings of the IEEE/CVF Conference on*
 612 *Computer Vision and Pattern Recognition*, pages 18720–
 613 18729, 2022.
- 614 [Thies *et al.*, 2016] Justus Thies, Michael Zollhofer, Marc
 615 Stamminger, Christian Theobalt, and Matthias Nießner.
 616 Face2face: Real-time face capture and reenactment of rgb
 617 videos. In *Proceedings of the IEEE conference on com-
 618 puter vision and pattern recognition*, pages 2387–2395,
 619 2016.
- 620 [Thies *et al.*, 2019] Justus Thies, Michael Zollhöfer, and
 621 Matthias Nießner. Deferred neural rendering: Image syn-
 622 thesis using neural textures. *Acm Transactions on Graph-
 623 ics (TOG)*, 38(4):1–12, 2019.
- 624 [Wang *et al.*, 2020] Run Wang, Felix Juefei-Xu, Qing Guo,
 625 Yihao Huang, Lei Ma, Yang Liu, and Lina Wang. Deep-
 626 tag: Robust image tagging for deepfake provenance. *arXiv*
 627 *preprint arXiv:2009.09869*, 3, 2020.
- 628 [Wang *et al.*, 2021] Run Wang, Felix Juefei-Xu, Meng Luo,
 629 Yang Liu, and Lina Wang. Faketagger: Robust safeguards
 630 against deepfake dissemination via provenance tracking.
- In *Proceedings of the 29th ACM International Conference* 631
 on *Multimedia*, pages 3546–3555, 2021. 632
- [Wu *et al.*, 2023] Xiaoshuai Wu, Xin Liao, and Bo Ou. 633
 Sepmark: Deep separable watermarking for unified 634
 source tracing and deepfake detection. *arXiv preprint* 635
arXiv:2305.06321, 2023. 636
- [Zhang *et al.*, 2018] Richard Zhang, Phillip Isola, Alexei A 637
 Efros, Eli Shechtman, and Oliver Wang. The unreasonable 638
 effectiveness of deep features as a perceptual metric. In 639
Proceedings of the IEEE conference on computer vision 640
and pattern recognition, pages 586–595, 2018. 641
- [Zhao *et al.*, 2021] Tianchen Zhao, Xiang Xu, Mingze Xu, 642
 Hui Ding, Yuanjun Xiong, and Wei Xia. Learning self- 643
 consistency for deepfake detection. In *Proceedings of the* 644
IEEE/CVF international conference on computer vision, 645
 pages 15023–15033, 2021. 646
- [Zhao *et al.*, 2023] Yuan Zhao, Bo Liu, Ming Ding, Baop- 647
 ing Liu, Tianqing Zhu, and Xin Yu. Proactive deepfake 648
 defence via identity watermarking. In *Proceedings of the* 649
IEEE/CVF winter conference on applications of computer 650
vision, pages 4602–4611, 2023. 651
- [Zheng *et al.*, 2021] Yinglin Zheng, Jianmin Bao, Dong 652
 Chen, Ming Zeng, and Fang Wen. Exploring temporal co- 653
 herence for more general video face forgery detection. In 654
Proceedings of the IEEE/CVF international conference on 655
computer vision, pages 15044–15054, 2021. 656
- [Zhu *et al.*, 2018] Jiren Zhu, Russell Kaplan, Justin Johnson, 657
 and Li Fei-Fei. Hidden: Hiding data with deep networks. 658
 In *Proceedings of the European conference on computer* 659
vision (ECCV), pages 657–672, 2018. 660
- [Zhuang *et al.*, 2022] Wanyi Zhuang, Qi Chu, Zhentao Tan, 661
 Qiankun Liu, Haojie Yuan, Changtao Miao, Zixiang Luo, 662
 and Nenghai Yu. Uia-vit: Unsupervised inconsistency- 663
 aware method based on vision transformer for face forgery 664
 detection. In *European Conference on Computer Vision*, 665
 pages 391–407. Springer, 2022. 666