

SpecGuard: Spectral Projection-based Advanced Invisible Watermarking

Anonymous ICCV submission

Paper ID 13406

Abstract

Watermarking embeds imperceptible patterns into images for authenticity verification. However, existing methods often lack robustness against various transformations primarily including distortions, image regeneration, and adversarial perturbation, creating real-world challenges. In this work, we introduce SpecGuard, a novel watermarking approach for robust and invisible image watermarking. Unlike prior approaches, we embed the message inside hidden convolution layers by converting from the spatial domain to the frequency domain using spectral projection of a higher frequency band that is decomposed by wavelet projection. Spectral projection employs Fast Fourier Transform approximation to transform spatial data into the frequency domain efficiently. In the encoding phase, a strength factor enhances resilience against diverse attacks, including adversarial, geometric, and regeneration-based distortions, ensuring the preservation of copyrighted information. Meanwhile, the decoder leverages Parseval's theorem to effectively learn and extract the watermark pattern, enabling accurate retrieval under challenging transformations. We evaluate the proposed SpecGuard based on the embedded watermark's invisibility, capacity, and robustness. Comprehensive experiments demonstrate the proposed SpecGuard outperforms the state-of-the-art models. To ensure reproducibility, we release the full code on [GitHub](#).

1. Introduction

With the rapid advancement of digital media and artificial intelligence, concerns regarding image authenticity, copyright protection, and content integrity have become more challenging than ever [7, 15, 31]. Moreover, the widespread availability of the latest image manipulation tools [3, 13] enables malicious tamperers to easily forge and redistribute digital content without authorization, posing a significant threat to ownership verification [28]. This growing risk emphasizes the need for reliable techniques for secure authentication and detection of unauthorized modifications of the original information.

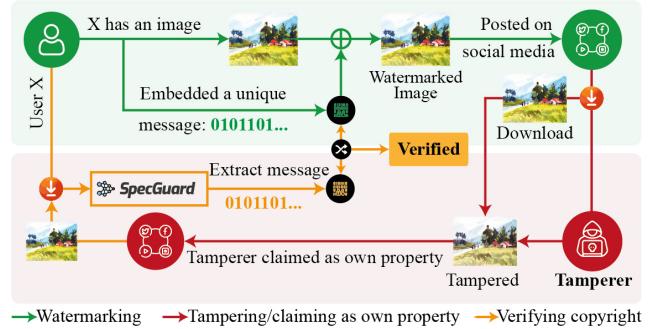


Figure 1. Image authentication using our proposed SpecGuard.

Recently, invisible watermarking has gained significant attention as a prominent defense mechanism for media authentication by embedding invisible messages into images to verify authenticity [34, 51]. In fact, invisible watermarks are preferred for preserving image quality and resisting tampering. These watermarks are unique to the creator and enable tamper verification by comparing the retrieved watermark to the original, as the high-level process is presented in Fig. 1. Traditional watermarking methods often rely on transformation techniques [17, 32]. Deep learning approaches like StegaStamp [44], Stable Signature [34], and HiDDeN [58] provide end-to-end solutions for message embedding. However, these methods often struggle with fragility in handling common image processing operations such as resizing, cropping, compression, and noise addition, which can distort or erase the embedded watermark. Additionally, the performance of watermark embedding and extraction often remains vulnerable to attacks with noise injection, blurring, contrasting, and rotation [6].

To address the aforementioned challenges, we introduce a novel robust, and invisible image watermarking method named SpecGuard. SpecGuard is designed to overcome the fundamental trade-offs [6] between imperceptibility, and robustness. Our proposed SpecGuard strategically embeds watermark information in the spectral domain, leveraging wavelet-based decomposition to distribute the watermark across high-frequency components. Unlike traditional frequency domain watermarking techniques [50, 53] that are

038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065

066 easily disrupted by common image manipulations, Spec-
 067 Guard maintains imperceptibility while significantly im-
 068 proving robustness against a wide range of transformations.

069 Overall, our proposed SpecGuard addresses the current
 070 limitations of the previous watermarking methods by pro-
 071 viding a robust, imperceptible watermarking technique that
 072 maintains integrity under diverse manipulations, signifi-
 073 cantly enhancing digital content security and authenticity
 074 verification. Our key contributions are as follows:

- 075 • We introduce a novel watermarking approach that em-
 076 beds message bits in high-frequency spectral components
 077 via wavelet and spectral projection inside hidden convolu-
 078 tional layers, ensuring robustness against various trans-
 079 formations and adversarial attacks.
- 080 • We adapt Parseval’s theorem [19] as a learnable threshold
 081 to optimize SpecGuard and spectral masking for robust
 082 watermark bit recovery under diverse transformations in-
 083 cluding distortions, regeneration, and adversarial attacks,
 084 proven through the experimental results.
- 085 • Our extensive evaluations demonstrate SpecGuard’s su-
 086 perior bit embedding capacity and producing better invis-
 087 ible watermarked images, surpassing the performance of
 088 state-of-the-art (SOTA) methods.

089 2. Related Works

090 Watermarking an image has been a widely researched topic
 091 for securing the ownership and verifying authenticity of digi-
 092 tal content [40]. Traditional watermarking techniques typi-
 093 cally embed invisible [45] or visible [5] watermarks into im-
 094 ages, which can later be extracted or detected to verify the
 095 content’s originality. These methods can be broadly classi-
 096 fied into spatial-domain [42, 47] and frequency-domain [11]
 097 watermarking, while some are based on combined meth-
 098 ods [41, 54]. However, researchers recently proposed many
 099 advanced models [22, 26, 30, 43] for effective watermark
 100 removal. To face this growing challenge, researchers intro-
 101 duced different methods [2, 21, 27, 55, 58] as alternatives to
 102 deep learning-based encoders or decoders to produce more
 103 robust image watermarking. Furthermore, iterative mod-
 104 els have demonstrated competitive performance [20, 33],
 105 particularly in robustness against a wide range of trans-
 106 formations. In addition, with the rise of generative meth-
 107 ods, researchers used the watermark-labeled data for train-
 108 ing to learn how to produce watermarks [9, 23]. Also,
 109 models that combine generative methods with watermarking
 110 techniques show promise in effective image watermarking
 111 [24, 29, 36]. However, such approaches face limitations
 112 such as increased computational complexity and longer pro-
 113 cessing times. These approaches are also more vulnerable
 114 to adversarial attacks that can target and distort the embed-
 115 ded watermark without altering the content visibly.

3. Proposed Method: SpecGuard

We introduce SpecGuard, as illustrated in Fig. 2, which involves two fundamental modules: an “Encoder” for embedding the watermark and a “Decoder” for accurately extracting the watermark detailed in the following sections.

3.1. Encoder

By targeting high-frequency components, the encoder integrates a binary message M into the cover image I . Using wavelet projection (WP) [32] and a Fast Fourier Transform (FFT)-based spectral projection (SP) [17] approximation, the message M is inserted into specific frequency bands, minimizing perceptual impact.

Wavelet Projection. We use a wavelet projection to capture frequency and spatial localization features that describe an image across different scales, as shown in Eq. (1):

$$W(a, b) = \frac{1}{\sqrt{|a|}} \int_{-\infty}^{\infty} f(x) \psi\left(\frac{x-b}{a}\right) dx, \quad (1)$$

where $a \in \mathbb{R} \setminus \{0\}$, and $b \in \mathbb{R}$ denote the scaling and translation parameters, respectively. Here, $\psi_{a,b}(x)$ represents a rescaled and translated form of the mother wavelet ψ , defined as follows:

$$\psi_{a,b}(x) = \psi\left(\frac{x-b}{a}\right) \cdot \frac{1}{\sqrt{|a|}}, \quad (2)$$

where $\frac{1}{\sqrt{|a|}}$ functions as a normalization factor, guaranteeing that the energy of the wavelet is invariant to the scaling parameter a . Minimal values of a compress the wavelet, enabling the inspection of high-frequency components, whereas greater values of a elongate the wavelet, promoting low-frequency analysis. Since each mother wavelet ψ is built with zero mean and finite energy [10], it guarantees to maintain stability as follows:

$$\int_{-\infty}^{\infty} \psi(x) dx = 0, \quad \int_{-\infty}^{\infty} |\psi(x)|^2 dx < \infty, \quad (3)$$

where the wavelet projection from Eq. (1) decomposes the input into orthogonal wavelet sets using discrete scales and translations. For 2D inputs, the scaled and translated basis elements [1] are defined for each coordinate pair (u, v) :

$$\begin{aligned} \mathbf{S}_{LL} &= \phi(u, v) = \phi(u)\phi(v), & \mathbf{S}_{LH} &= \psi_H(u, v) = \psi(u)\phi(v), \\ \mathbf{S}_{HL} &= \psi_V(u, v) = \phi(u)\psi(v), & \mathbf{S}_{HH} &= \psi_D(u, v) = \psi(u)\psi(v), \end{aligned} \quad (4)$$

where, H , V , and D represent the horizontal, vertical, and diagonal decomposition direction, respectively. To depict the image at different resolutions, we define scaling and wavelet functions at scale j as shown below:

$$\begin{aligned} \phi_{j,m,n}(u, v) &= 2^{j/2} \phi\left(u - \frac{m}{2^j}, v - \frac{n}{2^j}\right), \\ \psi_{j,m,n}^d(u, v) &= 2^{j/2} \psi^d\left(u - \frac{m}{2^j}, v - \frac{n}{2^j}\right), \end{aligned} \quad (5)$$

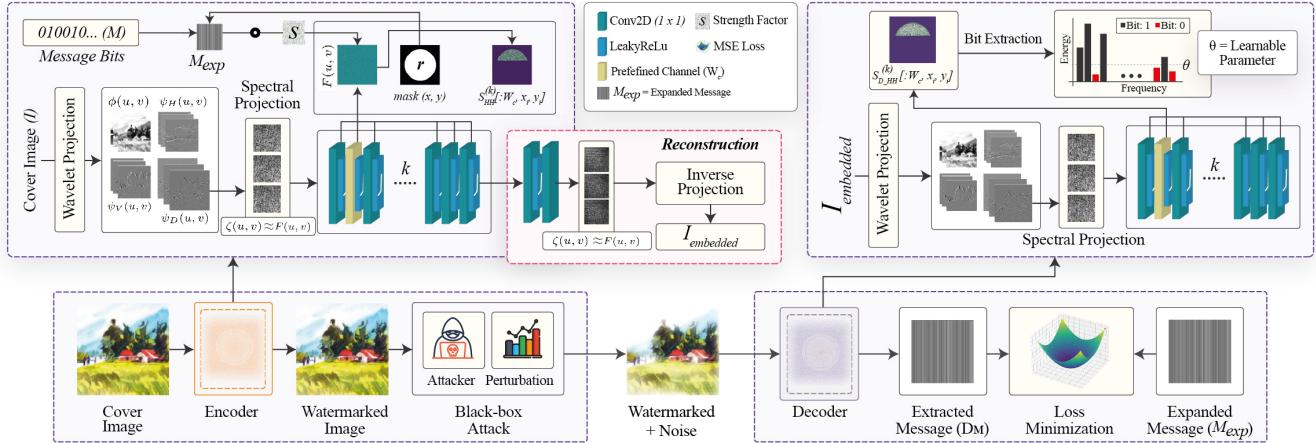


Figure 2. Architecture of the proposed SpecGuard watermarking method involves encoding a binary message M into the high-frequency band of the cover image I using wavelet and spectral projection and learning to decode the embedded message.

where $d \in \{H, V, D\}$ is the wavelet function direction that serves as discrete basis elements for multi-resolution analysis, capturing details across frequency bands and spatial locations. In Eq. (6), $T_{m,n}$ denotes the intensity or pixel value of the cover image I at spatial coordinates (m, n) . The discrete scaling function $W_\phi(j, u, v)$ (approximation at scale j) and the detail coefficients $W_\psi^d(j, u, v)$ for each direction are computed accordingly as follows:

$$\begin{aligned} W_\phi(j, u, v) &= \frac{1}{l} \sum_{m=0}^{l-1} \sum_{n=0}^{l-1} T_{m,n} \phi(m - u \cdot 2^{-j}, n - v \cdot 2^{-j}), \\ W_\psi^d(j, u, v) &= \frac{1}{l} \sum_{m=0}^{l-1} \sum_{n=0}^{l-1} T_{m,n} \psi^d(m - u \cdot 2^{-j}, n - v \cdot 2^{-j}), \end{aligned} \quad (6)$$

with l as the discrete region dimension, these coefficients capture multi-scale, multi-orientation image information, forming the basis of spectral features as follows:

$$\beta_j = \bigcup_{d \in \{H, V, D\}} (W_\phi(j, u, v) \cup W_\psi^d(j, u, v)). \quad (7)$$

This feature set β_j captures key frequency and spatial details across resolutions, forming the foundation for the watermark embedding process of our SpecGuard.

Selective Frequency Band Decomposition. To refine the embedding process, we segment the data into distinct frequency bands. The decomposition level κ is determined by the image complexity, calculated as follows:

$$\kappa = \lfloor \sqrt{\log(1 + N)} \rfloor, \quad (8)$$

where N denotes the total pixel count in the cover image I . And, each component β_j falls within a unique frequency band, yielding a total of $1 + 3\kappa$ distinct frequency bands as follows:

$$\beta_j = \phi_j(u, v) \cup \bigcup_{d \in \{H, V, D\}} \psi_j^d(u, v). \quad (9)$$

The components β_j , consisting of scaling functions $\phi_j(u, v)$ and wavelet functions $\psi_j^d(u, v)$, capture specific spatial frequency bands, enabling targeted high-frequency embedding. We translate the WP into disjoint intervals representing a unique frequency range to approximate the segmentation in the frequency domain:

$$\beta_j = \left\{ W_\psi^d(u, v) \mid u, v \in \left(\frac{j \cdot L}{\kappa}, \frac{(j+1) \cdot L}{\kappa} \right) \right\}, \quad (10)$$

where, L is the dimension of S_{HH} , and $W_\psi^d(u, v)$ represents wavelet values within segmented intervals. This frequency band partitioning mimics the frequency selectivity of wavelet sub-bands, enabling effective targeting of high-frequency regions for optimal embedding.

Approximation of Spectral Projection. We first apply spectral projection on the S_{HH} sub-band, transforming it into the spectral domain. Given a matrix $T(x, y)$ representing pixel intensities in S_{HH} , the spectral projection computes the spectral components $\zeta(u, v)$ as follows:

$$\zeta(u, v) = \frac{1}{L^2} \sum_x \sum_y T(x, y) \cdot \exp \left(-i \frac{2\pi}{L} (x \cdot u + y \cdot v) \right), \quad (11)$$

where L denotes the dimension of S_{HH} , $T(x, y)$ provides the intensity at each coordinate (x, y) which is equivalent to $W_\psi^d(u, v)$ in Eq. (6), i is the imaginary unit, and (u, v) are the spectral coordinates.

To approximate the spectral components using the FFT, we create a symmetrically extended version $\tilde{T}(x, y)$ of the original $N \times N$ matrix $T(x, y)$. This extension is achieved by mirroring $T(x, y)$ along its boundaries, doubling its size to $2N \times 2N$. Specifically, the original matrix occupies the top-left quadrant, with the remaining quadrants filled by reflecting $T(x, y)$ horizontally, vertically, and diagonally, respectively. This symmetric structure ensures that the FFT

212 yields only real values, allowing the spectral coefficients to
 213 be extracted directly from the real part of the FFT operation.
 214 Then, we apply the 2D FFT to $\tilde{T}(x, y)$ as follows:

$$F(u, v) = \frac{1}{(2N)^2} \sum_x \sum_y \tilde{T}(x, y) \cdot \exp\left(-i \frac{2\pi}{2N} (x \cdot u + y \cdot v)\right). \quad (12)$$

215 The SP coefficients are then approximated by taking the
 216 real part (Re) of F in the original $N \times N$ region as follows:
 217

$$\zeta(u, v) \approx \text{Re}(F(u, v)), \quad 0 \leq u, v < N. \quad (13)$$

218 Applying Eq. (13) to the sub-bands extracted from
 219 wavelet projection in Eq. (6), we achieve a computationally
 220 efficient spectral projection by leveraging the FFT approx-
 221 imation on a symmetrically extended matrix, maintaining
 222 effective embedding properties within the spectral domain.
 223

224 **SpecGuard Embedding Process.** The embedding process
 225 integrates the binary message M into the high-frequency
 226 band S_{HH} of the cover image I , enhancing robustness and
 227 imperceptibility through wavelet and spectral projection.
 228 Using the Eq. (6) and Eq. (13), the cover image I is de-
 229 composed into sub-bands S_{LL}, S_{LH}, S_{HL} , and S_{HH} within
 230 spectral domain, with S_{HH} providing high-frequency de-
 231 tails for embedding. A variable number k of convolutional
 232 layers with a $K \times K$ kernel, followed by LeakyReLU ac-
 233 tivation, are recursively applied to S_{HH} to refine spectral
 234 features as follows:

$$S_{HH}^{(n+1)} = \text{LeakyReLU}(\text{Conv}_{2D}(S_{HH}^{(n)}, K)), \quad n = 1, \dots, k. \quad (14)$$

235 The final output $S_{HH}^{(n+1)}$ from Eq. (14) represents the
 236 modified high-frequency band, primed for embedding.
 237

238 The message M , represented as a binary vector of length
 239 l ($M \in \{0, 1\}^l$), with batch size b and message length l ,
 240 is reshaped and expanded across channels c to align with
 241 $S_{HH}^{(n+1)}$. This transformation ensures M_{expanded} conforms to
 242 the dimension $[b, c, l]$, where each message is structured ac-
 243 cordingly.

244 To localize the embedding, we create a radial mask cen-
 245 tered at $(c_x, c_y) = (\frac{h}{2}, \frac{w}{2})$, where h and w represent the
 246 height and width of S_{HH} . The Euclidean distance $D(x_i, y_i)$
 247 from the center (c_x, c_y) is computed for each coefficient
 248 (x_i, y_i) . A binary mask is then generated within the pre-
 249 defined radius r based on the distance $D(x_i, y_i)$, such that
 250 if $D(x_i, y_i) \leq r$, the mask value is set to 1, allowing em-
 251 bedding in the corresponding region. Otherwise, the mask
 252 value is 0, restricting embedding to areas within a specified
 253 radius r , ensuring focus on high-frequency regions.

254 For each coordinate (x_i, y_i) where mask (x_i, y_i) is 1 and
 255 $W_c \in c$, the embedding operation is performed as follows:

$$S_{HH}^{(n+1)}[:, W_c, x_i, y_i] += M_{\text{expanded}}[:, W_c, i] \cdot s, \quad (15)$$

256 where s is the strength factor controlling embedding inten-
 257 sity and invisibility. After embedding, the modified coeffi-
 258 cients $S_{HH}^{(n+1)}$ undergo a final convolution and LeakyReLU
 259

260 using Eq. (14), by setting the value of $k = 1$ to harmonize
 261 the embedded message. Following this approach, Spec-
 262 Guard embeds the message into the spectral domain in a
 263 transformed form, differing from its original input represen-
 264 tation. By blending the message seamlessly into the spec-
 265 tral space based on the r , s , and W_c , it becomes inherently
 266 concealed within the domain, rendering its presence imper-
 267 ceptible. Without knowledge of r , s , and W_c , it becomes
 268 exceedingly challenging to localize the embedded message,
 269 further enhancing the security of the system. This trans-
 270 formation ensures the embedding process remains opaque
 271 to any adversarial attacker, effectively making SpecGuard a
 272 black-box system.

273 **Reconstruction.** SpecGuard encoder reconstructs the wa-
 274 termarked image I_{embedded} by inverse transformation restor-
 275 ing S_{HH} back into the spatial domain. The reconstruction
 276 process integrates the inverse wavelet projection (IWP) [32]
 277 and inverse spectral projection (ISP) [17], ensuring the em-
 278 bedded modifications are correctly translated into the spa-
 279 tial domain. To reconstruct the spatial domain image, S_{HH}
 280 is combined with the other sub-bands S_{LL}, S_{LH} , and S_{HL} .
 281 For the SP embedded in S_{HH} , the ISP is applied to recon-
 282 struct S_{HH} to spatial domain as follows:

$$S_{HH}(x, y) = \sum_{u=0}^{L-1} \sum_{v=0}^{L-1} \zeta(u, v) \cdot \exp\left(i \frac{2\pi}{L} (x \cdot u + y \cdot v)\right), \quad (16)$$

283 where $\zeta(u, v)$ represents spectral coefficients from the em-
 284 bedding process, L denotes the dimension of S_{HH} , and
 285 (x, y) are spatial coordinates. SpecGuard then reconstructs
 286 the watermarked image I_{embedded} using the IWP as follows:
 287

$$I_{\text{embedded}}(x, y) = \text{IWP}(S_{LL}, S_{LH}, S_{HL}, S_{HH}). \quad (17)$$

288 This process seamlessly embeds the watermark message
 289 M in the spectral domain, preserving the cover image I 's
 290 integrity. The inverse transformations that are expressed
 291 in Eq. (16) and Eq. (17) fully restore visual quality, main-
 292 taining all frequency components.
 293

3.2. Decoder

294 As shown in Algorithm 1, SpecGuard decoding process
 295 starts by applying wavelet projection (Eq. (1)) to the wa-
 296 termarked image I_{embedded} , separating it into low and high-
 297 frequency bands, where the high-frequency band $S_{D_{HH}}^{\text{high}}$
 298 contains the embedded message similar to the process in the
 299 encoding phase, particularly in Eq. (6). An approximation
 300 of the spectral projection using FFT as shown in Eq. (13) is
 301 then applied to $S_{D_{HH}}^{\text{high}}$ returning the transformed data $S_{D_{HH}}^{\text{sp}}$.
 302 Then, $S_{D_{HH}}$ is further refined through convolutional layers
 303 that captures the local features for message extraction.
 304

305 To extract the message, a radial mask is created to iso-
 306 late high-frequency areas within $S_{D_{HH}}$, targeting the em-
 307 bedded regions based on their distance from the center. The

Algorithm 1 SpecGuard decoder with wavelet, spectral projection with FFT approximation, and learnable threshold.

```

1: Input: Watermarked image  $I_{\text{embedded}}$ , learnable  $\theta$ , message length  $l$ ,  

   radius  $r$ , watermark channel  $W_c$ 
2: Output: Decoded binary message  $D_M$ 
3: Procedure: Apply Wavelet Projection on  $I_{\text{embedded}}$  to obtain  $S_{D_{LL}}$   

   (low-frequency) and  $S_{D_{HH}}^{\text{high}}$  (high-frequency)
4: Procedure: Spectral approximation with FFT ( $S_{D_{HH}}^{\text{high}}$ ):  

5: Separate even and odd indices:  $v = [x_{\text{even}}, \text{reverse}(x_{\text{odd}})]$ 
6: Compute FFT on  $v$ :  $V_{\text{complex}} = \text{FFT}(v)$ 
7:  $V_{\text{real}} = V_{\text{complex}} \cdot [\cos\left(\frac{-\pi k}{2N}\right), \sin\left(\frac{-\pi k}{2N}\right)]$  // Calculate Real
8:  $V_{\text{real}}[0] \leftarrow \frac{V_{\text{real}}[0]}{\sqrt{N \cdot 2}}, V_{\text{real}}[1:] \leftarrow \frac{V_{\text{real}}[1:]}{\sqrt{\frac{N}{2} \cdot 2}}$  // Energy preservation
9: Transpose result and repeat to obtain  $S_{D_{HH}}^{\text{sp}}$ 
10: Return  $S_{D_{HH}}^{\text{sp}}$ 
11: Procedure: Pass  $S_{D_{HH}}^{\text{sp}}$  through sequential layers as:  


$$S_{D_{HH}}^{(n+1)} = \text{LeakyReLU} \left( \text{Conv}_{2D} \left( S_{D_{HH}}^{sp(n)}, K \right) \right), n = 1, \dots, k,$$

12: Return  $S_{D_{HH}}^{(n+1)}$ 
13: Procedure: Extraction ( $S_{D_{HH}}^{(n+1)}$ ,  $l$ ):  

14: Set  $(c_x, c_y) = \left( \frac{H}{2}, \frac{W}{2} \right)$ 
15: Generate mask for high-frequency region within radius  $r$   

   for each coordinate  $(i, j)$  do:  

16:  $D(x_i, y_i) = \sqrt{(x_i - c_x)^2 + (y_i - c_y)^2}$  // Euclidian Distance
17: if  $D(x_i, y_i) \leq r$  then  

18:   Set mask[i, j] = 1
19: end if
20: end for
21: Extract mask:  $S_{D_{HH}}[:, W_c, \text{mask}[i, j]]$ 
22: Decode message using learnable  $\theta$ :  


$$D_M[i] = \begin{cases} 1 & \text{if Extracted}[i] > \theta \\ 0 & \text{otherwise} \end{cases}$$

23: Update  $\theta$  dynamically:  $\theta \leftarrow \theta - \eta \cdot \frac{\partial L_{\text{dec}}}{\partial \theta}$  // Optimizes robustness
24: Return  $D_M$ 

```

308 masked values are compared against a learnable threshold
 309 θ to decode each bit of the hidden message D_M . Here,
 310 θ serves as a threshold that adapts to the spectral patterns
 311 across the entire image, learning the distinct characteristics
 312 of the embedded watermark. From Parseval's theorem [19]
 313 ensures overall spectral and spatial energies remain equivalent,
 314 though local spectral energy distributions are altered
 315 by the watermark strength factor s .

316 The watermark's strength factor s ensures that the high-
 317 energy areas where the message M is embedded as "1"
 318 remain robust, experiencing a minimum distortion in such
 319 conditions. Moreover, this threshold can be optimized for
 320 better bit recovery accuracy during training. As θ learns,
 321 it recognizes that areas encoded as "1" carry higher en-
 322 ergy and impact due to the strength factor s of Eq. (15),
 323 while areas marked as "0", softened by the LeakyReLU's
 324 minimal negative slope, have a lower intensity. Such a dy-

namic approach enables θ to identify and protect the em-
 325 bedded message M even when external disturbances oc-
 326 cur, preserving the watermark's structure within the water-
 327 marked image I_{embedded} . And, θ effectively learns to distin-
 328 guish high-energy watermark regions. Therefore, the em-
 329 bedded message is more recoverable under diverse attacks,
 330 and SpecGuard's decoder ensures valid watermark bit ex-
 331 traction. Theoretical explanation of Parseval theorem's [19]
 332 impact on message extraction is in the Supplementary.
 333

3.3. Loss Calculation for SpecGuard

To achieve the training objective of robust and invisible wa-
 335 termark embedding, a composite loss function is defined
 336 with two terms: encoder loss L_{enc} as expressed in Eq. (18)
 337 and decoder loss L_{dec} as expressed in Eq. (19).
 338

$$\min_{\theta} \mathbb{E}_{(I, M) \sim D} L_{\text{enc}}(I, I_{\text{embedded}}) = \|E_{\theta}(I, M) - I\|^2, \quad (18)$$

$$\min_{\theta} \mathbb{E}_{(I, M) \sim D} L_{\text{dec}}(M, D_M) = \|D_{\theta}(I_{\text{embedded}}) - M\|^2, \quad (19)$$

where $E_{\theta}(I, M)$ denotes the encoder output, embedding the
 342 message M into the cover image I to produce I_{embedded} . By
 343 minimizing L_{enc} , the encoder learns to embed the wa-
 344 termark invisibly, preserving the fidelity of the cover im-
 345 age. $D_{\theta}(I_{\text{embedded}})$ denotes the decoder's output from the
 346 watermarked image I_{embedded} . Minimizing L_{dec} enables the
 347 decoder to reliably retrieve the embedded message under
 348 varying conditions, such as noise and transformation.
 349

The total loss L as shown in Eq. (20) used for optimizing
 350 the model combines these terms, balancing invisibility and
 351 robustness through weighted coefficients as follows:
 352

$$\min_{\theta} L = \lambda_{\text{enc}} L_{\text{enc}} + \lambda_{\text{dec}} L_{\text{dec}}, \quad (20)$$

where λ_{enc} and λ_{dec} control the relative importance of visual
 354 fidelity and message recoverability.
 355

4. Experimental Results**4.1. Dataset**

SpecGuard is trained on the MS-COCO dataset [25], which
 358 contains 25K images. To evaluate the robustness of the wa-
 359 termarking methods including our SpecGuard against dif-
 360 ferent types of attacks, such as distortions, regenerations,
 361 and adversarial attacks, we used three datasets: Diffu-
 362 sionDB [49], MS-COCO [25], and DALL-E3 ¹. Each of
 363 these datasets has a unique distribution of prompt words.
 364 We also ensured that no unethical or violent terms were in-
 365 cluded in the prompts. We randomly picked 200 images
 366 from MS-COCO [25] and applied watermark using Spec-
 367 Guard for further verifying the robustness after uploading
 368 on various social media platforms and applying AI-based
 369 Photoshop Neural Filters (PNFs) ². The PNFs include depth
 370

¹<https://huggingface.co/datasets/OpenDatasets/dalle-3-dataset>

²<https://www.adobe.com/products/photoshop/neural-filter.html>

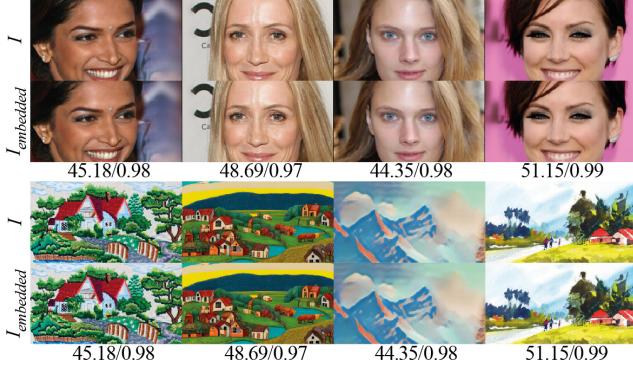


Figure 3. Some best results for cover vs watermarked images with PSNR/SSIM (\uparrow) scores showing minimal visual degradation when watermarked using proposed SpecGuard.

blur, artistic style transfer, super zoom, JPEG artifact reduction, and colorization. For the super zoom filter, we set the ‘Sharpen’ and ‘Noise Reduction’ parameters to 15. For all other filters, we used the default settings.

4.2. Implementation

We used CUDA v11.3 and PyTorch with a batch size of 32 and the Adam optimizer on a multiple NVIDIA RTX 2080-equipped server. Mean Squared Error (MSE) and Bit Recovery Accuracy (BRA) are used for loss and accuracy calculation. We used Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index (SSIM), Fréchet Inception Distance (FID), and MSE to evaluate perceptual quality. Our model is trained for 300 epochs, with the decoder learning rate set to 1×10^{-3} , reduced by half every 100 steps, and the encoder learning rate is set to 1×10^{-2} without scheduling. We set our watermark radius (r), strength factor (s), initial learning parameter (θ), and the number of convolutional layers (k) to 100, 20, 0.001, and 32, respectively. This setup is applied with a message bit length (BL) of 48, 64, 128, and 256. Initially, decoder loss weight (λ_{dec}) and encoder loss weight (λ_{enc}) are set to 1.0 and 0.7, respectively. For assessing the robustness of watermarking methods against diverse attacks, we inherited the experimental setups from Waves [6] and used effective metrics such as “Quality at 95% Performance (Q@0.95P)”, “Quality at 70% Performance (Q@0.7P)”, “Avg P” and “Avg Q.” Here, Q@0.95P and Q@0.7P indicate the level of image quality degradation required for watermark detection accuracy to reach 95% and 70%, respectively. The average performance (Avg P) metric represents the mean detection accuracy across various attack strengths, while the average quality degradation (Avg Q) measures the overall impact of attacks on image quality [6, 34].

4.3. Watermark Invisibility

To evaluate the invisibility of the embedded watermark, we conducted perceptual and quantitative assessments us-

| Metrics | 256 × 256 | | 512 × 512 | | 1024 × 1024 | |
|------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| | CelebA-HQ MS-COCO |
| PSNR \uparrow | 40.361 | 40.320 | 44.651 | 44.680 | 48.170 | 48.081 |
| SSIM \uparrow | 0.9889 | 0.9888 | 0.9927 | 0.9927 | 0.9937 | 0.9936 |
| FID \downarrow | 16.451 | 16.690 | 16.972 | 17.020 | 17.446 | 16.955 |
| MSE \downarrow | 0.0002 | 0.0002 | 0.0001 | 0.0001 | 0.0001 | 0.0001 |

Table 1. Perceptual quality evaluation for watermarked images using SpecGuard across resolutions and datasets.

| Methods | Venue | BL | PSNR \uparrow | SSIM \uparrow | FID \downarrow | BRA \uparrow |
|-----------------------|------------|-----|-----------------|-----------------|------------------|----------------|
| Tree-Ring [50] | NeurIPS'23 | 64 | 32.33 | 0.91 | 17.7 | 0.98 |
| | | 128 | 32.10 | 0.90 | 17.8 | 0.96 |
| | | 256 | 31.85 | 0.89 | 17.9 | 0.94 |
| | | 64 | 30.00 | 0.89 | 19.6 | 0.98 |
| Stable Signature [34] | ICCV'23 | 128 | 29.80 | 0.88 | 19.7 | 0.96 |
| | | 256 | 29.50 | 0.87 | 19.8 | 0.96 |
| Yang et al. [53] | CVPR'24 | 64 | 31.45 | 0.90 | 18.2 | 0.98 |
| | | 128 | 31.20 | 0.89 | 18.3 | 0.93 |
| | | 256 | 30.95 | 0.88 | 18.4 | 0.89 |
| SleeperMark [48] | CVPR'25 | 64 | 31.80 | 0.92 | 18.0 | 0.97 |
| | | 128 | 31.60 | 0.91 | 18.1 | 0.93 |
| | | 256 | 31.35 | 0.90 | 18.2 | 0.87 |
| HiDDeN [58] | ECCV'18 | 64 | 32.01 | 0.88 | 19.7 | 0.98 |
| | | 128 | 31.80 | 0.87 | 19.8 | 0.85 |
| | | 256 | 31.50 | 0.86 | 19.9 | 0.82 |
| StegaStamp [44] | CVPR'20 | 64 | 28.50 | 0.91 | 17.9 | 0.99 |
| | | 128 | 28.20 | 0.90 | 18.0 | 0.98 |
| | | 256 | 28.00 | 0.89 | 18.1 | 0.94 |
| MBRS [16] | ACM MM'21 | 64 | 38.20 | 0.96 | 17.9 | 0.98 |
| | | 128 | 37.90 | 0.95 | 18.0 | 0.96 |
| | | 256 | 37.50 | 0.94 | 18.2 | 0.94 |
| FIN [8] | AAAI'23 | 64 | 36.70 | 0.95 | 18.3 | 0.97 |
| | | 128 | 36.40 | 0.94 | 18.4 | 0.96 |
| | | 256 | 36.10 | 0.93 | 18.5 | 0.96 |
| MuST [46] | AAAI'24 | 64 | 41.20 | 0.97 | 17.5 | 0.98 |
| | | 128 | 40.90 | 0.96 | 17.6 | 0.93 |
| | | 256 | 40.50 | 0.95 | 17.8 | 0.90 |
| EditGuard [56] | CVPR'24 | 64 | 41.56 | 0.97 | 17.8 | 0.98 |
| | | 128 | 41.30 | 0.96 | 17.9 | 0.97 |
| | | 256 | 40.90 | 0.95 | 18.0 | 0.97 |
| SpecGuard | Ours | 64 | 42.59 | 0.98 | 17.2 | 0.99 |
| | | 128 | 42.89 | 0.99 | 17.0 | 0.99 |
| | | 256 | 40.86 | 0.99 | 17.6 | 0.98 |

*BL: Bit Length, BRA: Bit Recovery Accuracy

Table 2. Comparison of SOTA pre-processing and post-processing watermarking methods with SpecGuard without attacks.

ing SpecGuard. As shown in Fig. 3, there is no noticeable perceptual degradation between the cover and watermarked images, confirming that the watermark remains imperceptible to the human eye. For a more comprehensive evaluation, we created three subsets of different image sizes ranging between 256 to 1024 with images from the MS-COCO [25] and CelebA-HQ [18] datasets and applied the SpecGuard watermarking method to compare the average PSNR values between the cover and watermarked images, as in Tab. 1.

For quantitative evaluation, we further compare the performance of SpecGuard with the SOTA pre-processing and post-processing watermarking methods. As presented in Tab. 2, SpecGuard achieves the highest PSNR of 42.89 when the bit length was 128. Additionally, it attains the highest SSIM of 0.99 at a BL of 128 and 256 among all compared methods, indicating minimal visual distortion. Additionally, SpecGuard achieved the lowest FID of 17.0 and the highest BRA of 0.99, ensuring strong

| Attack Type | Tree-Ring [50] | | | | Stable Signature [34] | | | | StegoStamp [44] | | | | SpecGuard (Ours) | | | | | |
|--------------|-----------------|--------|-------|-------|-----------------------|--------|-------|-------|-----------------|--------|-------|-------|------------------|--------|-------|-------|-------|-------|
| | Q@0.95P | Q@0.7P | Avg P | Avg Q | Q@0.95P | Q@0.7P | Avg P | Avg Q | Q@0.95P | Q@0.7P | Avg P | Avg Q | Q@0.95P | Q@0.7P | Avg P | Avg Q | | |
| Distortions | Rotation | 0.464 | 0.521 | 0.375 | 0.648 | 0.624 | 0.702 | 0.594 | 0.650 | 0.423 | 0.498 | 0.357 | 0.616 | 0.863 | 0.863 | 0.687 | 0.653 | |
| | Crop | 0.592 | 0.592 | 0.332 | 0.463 | inf | inf | 0.995 | 0.461 | 0.602 | 0.602 | 0.540 | 0.451 | 0.812 | 0.812 | 0.998 | 0.742 | |
| | Bright | inf | inf | inf | 0.304 | inf | inf | 0.998 | 0.305 | inf | inf | 0.998 | 0.317 | inf | inf | 0.998 | 0.466 | |
| | Contrast | inf | inf | 0.998 | 0.243 | inf | inf | 0.998 | 0.243 | inf | inf | 0.998 | 0.231 | inf | inf | 0.998 | 0.556 | |
| | Blur | 0.861 | 1.112 | 0.563 | 1.221 | — inf | — inf | 0.000 | 1.204 | 0.848 | 0.962 | 0.414 | 1.000 | 0.921 | inf | 1.000 | 1.452 | |
| | Noise | 0.548 | inf | 0.980 | 0.395 | 0.402 | 0.520 | 0.870 | 0.390 | inf | inf | 1.000 | 0.360 | inf | inf | 0.999 | 0.568 | |
| | JPEG | 0.499 | 0.499 | 0.929 | 0.284 | 0.485 | 0.485 | 0.793 | 0.284 | inf | inf | 0.998 | 0.263 | inf | inf | 1.000 | 0.495 | |
| | Geo | 0.525 | 0.593 | 0.277 | 0.768 | 0.850 | inf | 0.937 | 0.767 | 0.663 | 0.693 | 0.396 | 0.733 | 0.869 | 0.869 | 0.865 | 0.623 | |
| | Deg | 0.620 | inf | 0.892 | 0.694 | 0.206 | 0.369 | 0.300 | 0.679 | 0.826 | 0.975 | 0.852 | 0.664 | 0.895 | 1.141 | 0.915 | 0.749 | |
| | Combine | 0.539 | 0.751 | 0.403 | 0.908 | 0.538 | 0.691 | 0.334 | 0.900 | 0.945 | 1.101 | 0.795 | 0.870 | 0.979 | 1.256 | 0.911 | 0.952 | |
| Regeneration | Regen-Diff | — inf | 0.307 | 0.612 | 0.323 | — inf | — inf | 0.001 | 0.300 | 0.331 | inf | 0.943 | 0.327 | inf | inf | 0.982 | 0.477 | |
| | Regen-DiffP | inf | 0.307 | 0.601 | 0.327 | — inf | — inf | 0.001 | 0.303 | 0.333 | inf | 0.940 | 0.329 | inf | inf | 0.982 | 0.562 | |
| | Regen-VAE | 0.578 | 0.578 | 0.832 | 0.348 | 0.545 | 0.545 | 0.516 | 0.339 | inf | inf | 1.000 | 0.343 | inf | inf | 0.995 | 0.521 | |
| | Regen-KLVAE | inf | inf | 0.990 | 0.233 | 6 | — inf | 0.176 | 0.217 | 0.206 | inf | inf | 1.000 | 0.240 | inf | inf | 0.990 | 0.492 |
| | Rinse-2xDiff | — inf | 0.333 | 0.510 | 0.357 | — inf | — inf | 0.001 | 0.332 | 0.391 | inf | 0.941 | 0.366 | inf | inf | 0.993 | 0.561 | |
| | Rinse-4xDiff | — inf | 0.355 | 0.443 | 0.466 | — inf | — inf | 0.000 | 0.438 | 0.388 | inf | 0.909 | 0.477 | inf | inf | 0.992 | 0.533 | |
| Adversarial | AdvEmbG-KLVAE8 | — inf | 0.164 | 0.448 | 0.253 | inf | inf | 0.998 | 0.249 | inf | inf | 1.000 | 0.232 | inf | inf | 1.000 | 0.456 | |
| | AdvEmbB-RN18 | 0.241 | inf | 0.953 | 0.218 | inf | inf | 0.999 | 0.212 | inf | inf | 1.000 | 0.196 | inf | inf | 1.000 | 0.467 | |
| | AdvEmbB-CLIP | 0.541 | inf | 0.932 | 0.549 | inf | inf | 0.999 | 0.541 | inf | inf | 1.000 | 0.488 | inf | inf | 1.000 | 0.436 | |
| | AdvEmbB-KLVAE16 | 0.195 | inf | 0.888 | 0.238 | inf | inf | 0.997 | 0.233 | inf | inf | 1.000 | 0.206 | inf | inf | 1.000 | 0.482 | |
| | AdvEmbB-SdxIVAE | 0.222 | inf | 0.934 | 0.221 | inf | inf | 0.998 | 0.219 | inf | inf | 1.000 | 0.204 | inf | inf | 1.000 | 0.492 | |
| | AdvCls-UnWM&WM | — inf | 0.102 | 0.499 | 0.145 | inf | inf | 0.999 | 0.101 | inf | inf | 1.000 | 0.101 | inf | inf | 1.000 | 0.497 | |
| | AdvCls-Real&WM | inf | inf | 1.000 | 0.047 | inf | inf | 0.998 | 0.092 | inf | inf | 1.000 | 0.106 | inf | inf | 1.000 | 0.427 | |
| | AdvCls-WM1&WM2 | — inf | 0.101 | 0.492 | 0.139 | inf | inf | 0.999 | 0.084 | inf | inf | 1.000 | 0.129 | inf | inf | 1.000 | 0.441 | |

Table 3. Robustness comparison various across attacks using Q@0.95P(↑), Q@0.7P(↑), Avg P(↑) and Avg Q(↑). Here, ‘inf’ denotes that no attack was sufficient to degrade performance below the threshold, indicating strong robustness, whereas ‘-inf’ signifies that even the weakest attack caused detection to fall below the threshold, reflecting weak robustness.

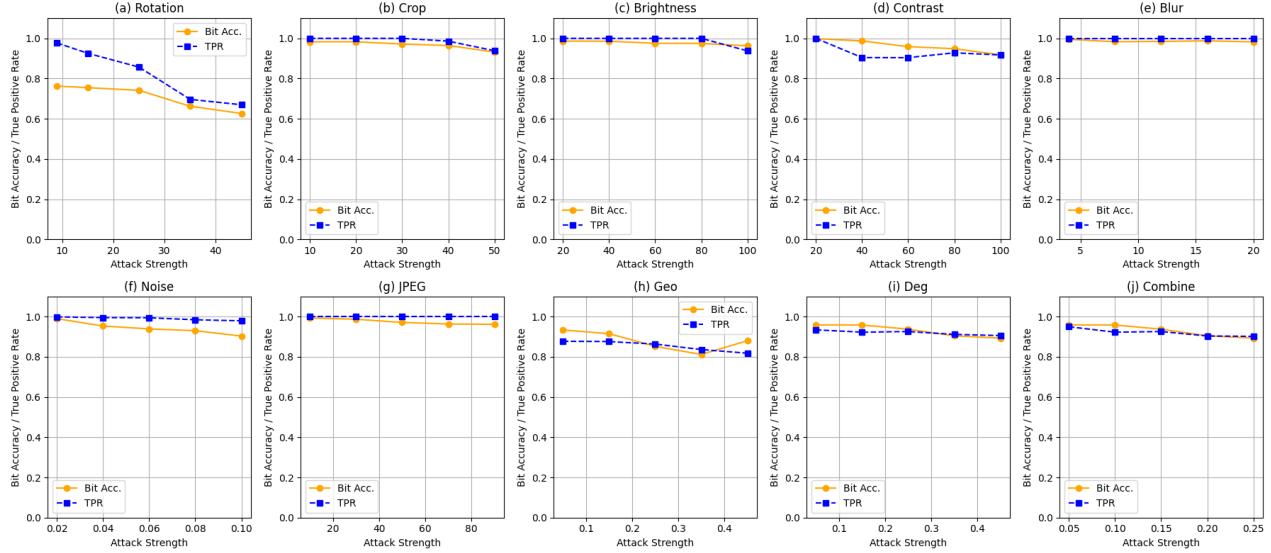


Figure 4. Robustness validation of our proposed SpecGuard under different distortion attacks, including geometric transformations: Geo (rotation, cropping), photometric modifications (brightness, contrast), and degradations: Deg (blur, noise, JPEG compression).

robustness while maintaining imperceptibility. Overall, our results demonstrate that SpecGuard outperforms both pre-processing and post-processing watermarking methods, achieving superior imperceptibility and robustness.

4.4. Capacity

To evaluate embedding capacity, we examined SpecGuard across different bit lengths and compared it with SOTA wa-

termarking methods. Our experiments with 64, 128, and 256 bits demonstrate SpecGuard’s high capacity of bit embedding while maintaining perceptual quality and robustness, as shown in Tab. 2. Notably, it achieves a PSNR of 42.89, the highest among all methods, along with the highest BRA of 0.99 and the lowest FID of 17.0 at 128 bits, ensuring minimal visual impact. This adaptability to different bit lengths without quality loss makes SpecGuard ideal

| Modules | PSNR/SSIM↑ | BRA↑ | Modules | PSNR/SSIM↑ | BRA↑ |
|--------------|-------------|------|------------------------------|-------------------|-------------|
| WP(L_1) | 40.51/0.96 | 0.92 | WP(L_1)+SP _{FA} | 42.89/0.99 | 0.99 |
| WP(L_2) | 38.15/0.93 | 0.87 | WP(L_2)+SP _{FA} | 36.25/0.92 | 0.89 |
| Attacks | PSNR/SSIM↑ | BRA↑ | Attacks | PSNR/SSIM↑ | BRA↑ |
| Rotate (45°) | 12.15/21.31 | 0.82 | Rotate (90°) | 11.15/19.31 | 0.65 |
| Blur (0.3) | 35.01/0.95 | 0.98 | Blur (0.6) | 30.11/0.91 | 0.98 |
| Geo (0.3) | 12.08/0.50 | 0.93 | Geo (0.6) | 10.25/0.45 | 0.86 |

*WP: Wavelet Projection, SP: Spectral Projection, FA: FFT Approximation

Table 4. Ablation studies on the proposed SpecGuard for across various configurations, setting $M = 128$, $r = 100$, and $s = 20$.

| Platform | PSNR/SSIM↑ | BRA↑ | PS Filters | PSNR/SSIM↑ | BRA↑ |
|-------------|-------------------|-------------|-------------------|-------------------|-------------|
| Facebook | 48.56/0.97 | 0.97 | Depth Blur | 25.25/0.89 | 0.85 |
| LinkedIn | 47.55/0.97 | 0.96 | StyleT. | 25.12/0.84 | 0.85 |
| Instagram | 48.56/0.98 | 0.98 | Super Zoom | 36.15/0.88 | 0.95 |
| WhatsApp | 42.10/0.96 | 0.97 | JPEG Artifacts | 31.01/0.85 | 0.94 |
| X (Twitter) | 49.25/1.00 | 0.99 | Colorize | 23.15/0.82 | 0.92 |

Table 5. Evaluation of SpecGuard’s robustness across Photoshop filters and while uploaded on different social media platforms.

for applications requiring flexible watermark sizes. Unlike StegaStamp and HiDDeN, which suffer reduced BRA for higher message bits, SpecGuard consistently extracts bits across all tested lengths. SpecGuard’s theoretical watermark capacity is provided in the Supplementary.

4.5. Robustness

We evaluate watermarking robustness by analyzing detection performance against a range of diverse and challenging real-world attacks. Results demonstrate the strong robustness of SpecGuard across various attacks. For example, as presented in Tab. 3, against geometric distortions such as cropping and rotation, SpecGuard achieved an Avg P of 0.998 and 0.687, respectively. Similarly, across the combined distortion-based attacks, SpecGuard achieves an overall Avg P of 0.911 and Avg Q of 0.952, ensuring minimal quality loss while maintaining high detection accuracy. Notably, the high values of Q@0.95P and Q@0.7P indicate that SpecGuard can sustain reliable detection at strict performance thresholds, even under aggressive perturbations. Unlike prior methods that struggle with extreme transformations, SpecGuard shows remarkable robustness against regeneration-based attacks like Rinse-2xDiff [4] (an image is noised then denoised by Stable Diffusion v1.4 two times with strength as a number of timesteps, 20-100) and Regen-VAE [4], maintaining high Avg P. Similarly, under adversarial attacks, SpecGuard consistently secures watermark detectability, outperforming existing techniques across all tested scenarios. These results establish SpecGuard as a highly robust watermarking approach capable of preserving image integrity even under severe distortions and adversarial manipulations, ensuring watermark reliability across diverse attack types. More details about how the attacks are performed are provided in supplementary material. Further, our results in Fig. 4 highlight the strong robustness of Spec-

Guard against various distortion attacks compared to other SOTA watermarking methods.

Social Platforms and Photoshop Filters. SpecGuard’s robustness when images are shared across social media platforms and subjected to common Photoshop Neural Filters (PNFs) is shown in Tab. 5. SpecGuard consistently maintains high PSNR and SSIM values, with BRA values close to 0.99 on platforms such as X (formally Twitter), Instagram, and Facebook. Also, it shown strong resilience to various PNFs, such as Super Zoom and JPEG Artifacts achieving BRA of 0.95 and 0.94. The PSNR, SSIM, and BRA values are expected to decrease with the severity of image manipulation, as increased manipulation leads to loss of image authenticity. For example, as we applied 60% style transfer the PSNR and BRA decreased to 25.12 and 0.85. Similarly, the depth blur which excessively reduces the image clarity also causes the decrease of BRA to 0.85.

4.6. Ablation Study

We examined the impact of wavelet projection (WP) at different levels (L_1 and L_2) and its combination with spectral projection (SP) using FFT approximation (FA) in Tab. 4. As observed, the WP(L_1) + SP_{FA} configuration achieved the highest PSNR and SSIM values of 42.89 and 0.99, respectively, and BRA of 0.99, indicating improved watermark invisibility and robustness. In contrast, using WP alone at either L_1 or L_2 resulted in lower BRA, with values of 0.92 and 0.87, respectively, demonstrating that the combined WP + SP_{FA} approach significantly enhances performance. We also evaluated the robustness of SpecGuard under strong adversarial attacks identified in Tab. 4, such as rotation, blur, and geometric transformations. The results indicate that higher levels of attack severity, such as 90° rotation, lead to a more significant drop in PSNR, SSIM, and BRA, with values dropping to 11.15, 19.31, and 0.65, respectively. Despite this, the model shows relatively high resilience under moderate attack intensities, such as 45° rotation and low levels of blur and geometric distortion, achieving BRA values as high as 0.93 under geometric transformations at the 0.3 thresholds. More ablations on SpecGuard are in the supplementary.

5. Conclusion

We propose SpecGuard, a novel invisible watermarking method that ensures secure and robust information concealment. Unlike traditional approaches, SpecGuard remains highly resilient against diverse distortions, adversarial attacks, and regeneration-based transformations. Experimental results demonstrate its superior bit recovery accuracy of 99% maintaining high PSNR. By outperforming SOTA watermarking methods in both detection reliability and imperceptibility, SpecGuard establishes a new benchmark for watermarking under real-world constraints.

SpecGuard: Spectral Projection-based Advanced Invisible Watermarking

Supplementary Material

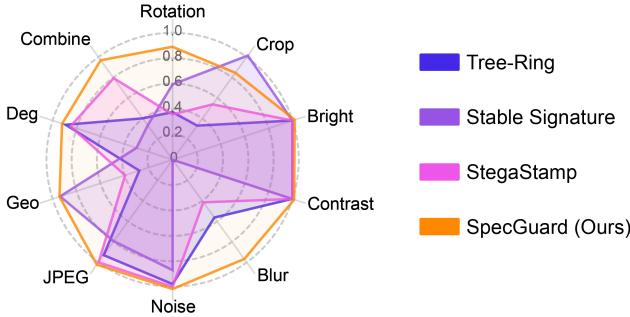


Figure 5. Comparison of SOTA watermarking methods in terms of average TPR@0.1%FPR (90% of watermarked images are correctly detected at 0.1% false positive rate) under different attacks.

525

6. Summary of Notations

526
527
528
529

To ensure clarity in understanding SpecGuard’s mathematical formulation, we summarize the key notations used throughout the methodology (Sec. 3) of the main paper. The complete set of notations is presented in Tab. 6.

530
531

7. Impact of Parseval’s Theorem in Message Extraction

532
533
534
535
536
537

To achieve robust and efficient decoding as detailed in Sec. 3.2 of the main paper, SpecGuard leverages Parseval’s theorem [19], a fundamental principle in signal processing, which establishes energy equivalence between spatial and spectral domains. Formally, Parseval’s theorem is defined as follows:

538

$$\sum_{x,y} |I(x,y)|^2 = \sum_{u,v} |\zeta(u,v)|^2, \quad (21)$$

539
540
541

where $I(x,y)$ denotes spatial-domain pixel intensities, and $\zeta(u,v)$ represent their corresponding spectral-domain coefficients.

542
543
544
545

In SpecGuard, watermark embedding modifies selected spectral coefficients, introducing subtle local energy variations. The embedding process employs a strength factor s , adjusting spectral energy differences as follows:

546

$$\zeta_{\text{embedded}}(u,v) = \zeta(u,v) + s \cdot W(u,v), \quad (22)$$

547
548
549
550
551

where $\zeta_{\text{embedded}}(u,v)$ denotes modified coefficients and $W(u,v)$ is the spectral-domain watermark signal. Although local energy distribution is altered, the overall signal energy remains constant as guaranteed by Parseval’s theorem as follows:

552

$$\sum_{x,y} |I(x,y)|^2 = \sum_{u,v} |\zeta_{\text{embedded}}(u,v)|^2. \quad (23)$$

During decoding, these local spectral energy variations, preserved due to total energy constancy, allow stable watermark extraction. Specifically, the decoder computes spectral projections via FFT approximation to isolate embedded spectral energy patterns as follows:

$$S_{D_{HH}}^{\text{sp}} = \text{SpectralProjectionFFT}(S_{D_{HH}}^{\text{high}}). \quad (24)$$

The decoder subsequently employs a dynamically optimized threshold θ to differentiate watermark signals from noise as follows:

$$D_M[i] = \begin{cases} 1 & \text{if } S_{D_{HH}}^{\text{sp}}[i] > \theta, \\ 0 & \text{otherwise.} \end{cases} \quad (25)$$

The adaptive threshold θ is optimized via gradient descent during training, adapting to spectral energy distributions as follows:

$$\theta \leftarrow \theta - \eta \cdot \frac{\partial L_{\text{dec}}}{\partial \theta}, \quad (26)$$

where L_{dec} is the decoding loss, and η is the learning rate. Thus, Parseval’s theorem critically supports SpecGuard by preserving total spectral energy, enabling stable differentiation of watermark bits and reliable decoding even under diverse real-world image distortions and adversarial attacks.

8. Mathematical Proof

8.1. Proof for S_{HH} Band of Wavelet Projection.

Here we presented a proof of one of the wavelet projections S_{HH} from Eq. (4) based on the Eq. (6) of the main paper.

$$\psi_j^D(u) = 2^{j/2} \psi^D(2^j u), \quad //1D \text{ wavelet} \quad (576)$$

$$\psi_{j,m}^D(u) = 2^{j/2} \psi^D(2^j u - m), \quad //\text{Translation} \quad (578)$$

$$\psi_{j,m,n}^D(u, v) = 2^{j/2} \psi^D(2^j u - m) \cdot \psi^D(2^j v - n), \quad //2D \text{ wavelet} \quad (580)$$

$$S_{HH}(j, m, n) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I(u, v) \cdot \psi_{j,m}^D(u) \psi_{j,m,n}^D(v) du dv, \quad (582)$$

$$\psi_{j,m,n}^D(u, v) du dv, \quad //\text{Projection} \quad (583)$$

$$S_{HH}(j, m, n) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I(u, v) \cdot \left[2^{j/2} \psi^D(2^j u - m) \cdot \psi^D(2^j v - n) \right] du dv, \quad //\text{Substitution} \quad (585)$$

| Notation | Description |
|--|--|
| I | Cover image |
| I_{embedded} | Watermarked image |
| M | Watermark message |
| c | Number of channels (e.g., RGB has $c = 3$) |
| H, W | Height and width of the image |
| $W(a, b)$ | Wavelet transform of signal $f(x)$ |
| a, b | Scaling and translation parameters in wavelet transform |
| ψ | Mother wavelet function |
| d | Direction of each wavelet components derived from ψ |
| $\phi(u, v), \psi_H(u, v), \psi_V(u, v), \psi_D(u, v)$ | Every directional scaling and wavelet basis components |
| $S_{LL}, S_{LH}, S_{HL}, S_{HH}$ | Wavelet sub-bands (low and high frequency components) |
| β_j | Feature set capturing frequency and spatial details |
| κ | Decomposition level determined by image complexity |
| $T(x, y)$ | Pixel intensity in high-frequency sub-band S_{HH} |
| $\zeta(u, v)$ | Spectral projection coefficients |
| s | Strength factor controlling embedding intensity |
| (c_x, c_y) | Center coordinates of the image |
| $D(x_i, y_i)$ | Euclidean distance from the center |
| r | Radius of embedding region |
| W_c | Selected watermark channel for embedding |
| θ | Learnable threshold for watermark extraction |
| $F(u, v)$ | 2D Fast Fourier Transform (FFT) of the extended signal |
| $L_{\text{enc}}, L_{\text{dec}}$ | Encoder and decoder loss functions |

Table 6. Description of the notations we used in the Sec. 3 (main paper) to describe our proposed SpecGuard.

587

588
$$S_{HH}(j, m, n) = \sum_{p=0}^{l-1} \sum_{q=0}^{l-1} T_{m,n} \cdot \psi^D(2^j u - m) \\ \cdot \psi^D(2^j v - n), \quad //\text{Discretization}$$

589
590
591
592
$$W_{\psi}^d(j, u, v) = \frac{1}{l} \sum_{m=0}^{l-1} \sum_{n=0}^{l-1} T_{m,n} \cdot \psi^D(m - u \cdot 2^{-j}, \\ n - v \cdot 2^{-j}), \quad //\text{Normalized}$$

593 8.2. Maximum Theoretical Watermark Capacity

594 To determine the maximum theoretical watermark capacity of SpecGuard, we analyze the SpecGuard’s embedding
595 pipeline, which integrates wavelet projection and spectral
596 projection. The capacity derivation considers three key
597 stages: ‘wavelet projection,’ ‘spectral projection,’ and ‘wa-
598 termark distribution,’ with each stage affecting the number
599 of available coefficients for embedding.

600 **Impact of Wavelet Projection.** SpecGuard applies wavelet
601 projection at decomposition level L , dividing the image

| Activation Function | Radius (r) | PSNR↑ | SSIM↑ | BRA↑ |
|---------------------|----------------|--------------|-------------|-------------|
| ReLU | $r(50)$ | 39.54 | 0.93 | 0.97 |
| | $r(75)$ | 38.64 | 0.91 | 0.93 |
| | $r(100)$ | 37.96 | 0.91 | 0.95 |
| Tanh | $r(50)$ | 37.18 | 0.89 | 0.82 |
| | $r(75)$ | 35.33 | 0.85 | 0.78 |
| | $r(100)$ | 37.66 | 0.90 | 0.80 |
| LeakyReLU | $r(50)$ | 39.77 | 0.96 | 0.98 |
| | $r(75)$ | 40.28 | 0.97 | 0.98 |
| | $r(100)$ | 42.89 | 0.99 | 0.99 |

Table 7. Performance evaluation of SpecGuard for different radius size and activation functions while the Strength Factor is 20.

603 into sub-bands. The watermark is embedded in the high-
604 frequency sub-band, which retains fine image details and
605 ensures robustness against low-frequency distortions. The
606 spatial dimensions of the wavelet sub-band are reduced by
607 a factor of 2^L along both height and width, resulting in a
608 down-sampling effect.

609 The number of available coefficients after wavelet de-

| Activation Function | Strength Factor (s) | PSNR↑ | SSIM↑ | BRA↑ |
|---------------------|---------------------------|--------------|-------------|-------------|
| LeakyReLU | $s(5)$ | 40.79 | 0.98 | 0.97 |
| LeakyReLU | $s(10)$ | 39.51 | 0.96 | 0.97 |
| LeakyReLU | $s(15)$ | 38.14 | 0.95 | 0.99 |
| LeakyReLU | $s(20)$ | 42.89 | 0.99 | 0.99 |

Table 8. Impact of Strength Factor for the best combination of the activation function (LeakyReLU) and radius $r(100)$.

composition is as follows:

$$N_{WP} = \frac{H \times W}{4^L}, \quad (27)$$

where H and W are the image height and width, respectively. Including all image channels c , the total number of wavelet coefficients available for embedding is as follows:

$$N_{WP,\text{total}} = \frac{H \times W \times c}{4^L}. \quad (28)$$

Thus, increasing the decomposition level L reduces the available spatial coefficients exponentially, limiting embedding capacity.

Impact of Spectral Project. SpecGuard employs spectral projection using FFT to distribute the watermark in the spectral domain. The spectral coefficients are selectively utilized based on an adaptive mask that prioritizes mid-to-high-frequency components while avoiding low frequencies (which contain most perceptual information) and extremely high frequencies (which are prone to compression loss).

The fraction of spectral coefficients selected for watermarking is denoted as f_{spectral} where spectral coefficients are used in between 20% and 50% as follows:

$$0.2 \leq f_{\text{spectral}} \leq 0.5. \quad (29)$$

After spectral projection following Eq. (28), the number of coefficients available for embedding is as follows:

$$N_{SP} = f_{\text{spectral}} \times N_{WP,\text{total}} = f_{\text{spectral}} \times \frac{H \times W \times c}{4^L}. \quad (30)$$

A higher f_{spectral} increases embedding capacity but may reduce robustness to compression and noise, while a lower f_{spectral} focuses on the most resilient coefficients but limits capacity.

Watermark Distribution and Final Capacity. The watermark is distributed across the selected spectral coefficients f_{spectral} using a weighting scheme, where each coefficient can embed multiple bits. Let N_b represent the number of watermark bits per selected coefficient f_{spectral} . The total embedded bits are then as follows:

$$C_{\text{total}} = N_b \times N_{SP}. \quad (31)$$

Substituting N_{SP} , the final maximum theoretical watermark capacity of SpecGuard is as follows:

$$C_{\max}(H, W, c, L, f_{\text{spectral}}, N_b) = \frac{H \times W \times c}{4^L} \times f_{\text{spectral}} \times N_b. \quad (32)$$

The watermark capacity scales proportionally with the image dimensions $H \times W$ and the number of channels c , ensuring that larger images provide greater embedding space. However, higher wavelet decomposition levels L reduce the available capacity exponentially due to the 4^L down-sampling effect. The fraction of spectral coefficients selected for embedding, denoted as f_{spectral} , controls how much of the frequency domain is utilized, balancing capacity and robustness. Additionally, the bit depth N_b determines the number of bits embedded per coefficient, directly influencing the total watermark payload.

Thus, SpecGuard achieves a flexible balance between capacity and robustness by leveraging adaptive spectral selection and wavelet decomposition, ensuring resilience under various transformations and attacks.

9. Impact of Hyperparamters

The performance of SpecGuard is influenced by several key hyperparameters, including the activation function, radius size (r), and strength factor (s). Each parameter plays a vital role in balancing the trade-off between perceptual quality, robustness, and watermark recovery accuracy. In addition to the ablation studies shown in Section 4.5 in the main paper, here we analyze the effect of the hyperparameters individually by conducting experiments under controlled conditions and report the findings in Tab. 7 and Tab. 8. All the experiments presented here were conducted using a 128-bit watermark message.

9.1. Activation Function and Radius

Table 7 highlights the performance of SpecGuard with various activation functions, including ReLU [12], Tanh [38], and LeakyReLU [52], while keeping the strength factor s fixed at 20. Among these, LeakyReLU outperforms others in terms of PSNR, SSIM, and bit recovery accuracy values across different radius sizes. Notably, with a radius r of 100, LeakyReLU achieves a PSNR and SSIM of 42.89 and 0.99, respectively, with a bit recovery accuracy of 0.99. Overall, the results indicate the effectiveness of LeakyReLU for robust and invisible watermarking compared to ReLU and Tanh. While testing with different r , such as 50 and 75, we observed a slightly lower perceptual quality and bit recovery accuracy. Therefore, we propose the SpecGuard with a combination of LeakyReLU, r of 100 and s of 20.

9.2. Strength Factor

Table 8 investigates the impact of the strength factor (s) using the best combination of LeakyReLU and radius $r(100)$.

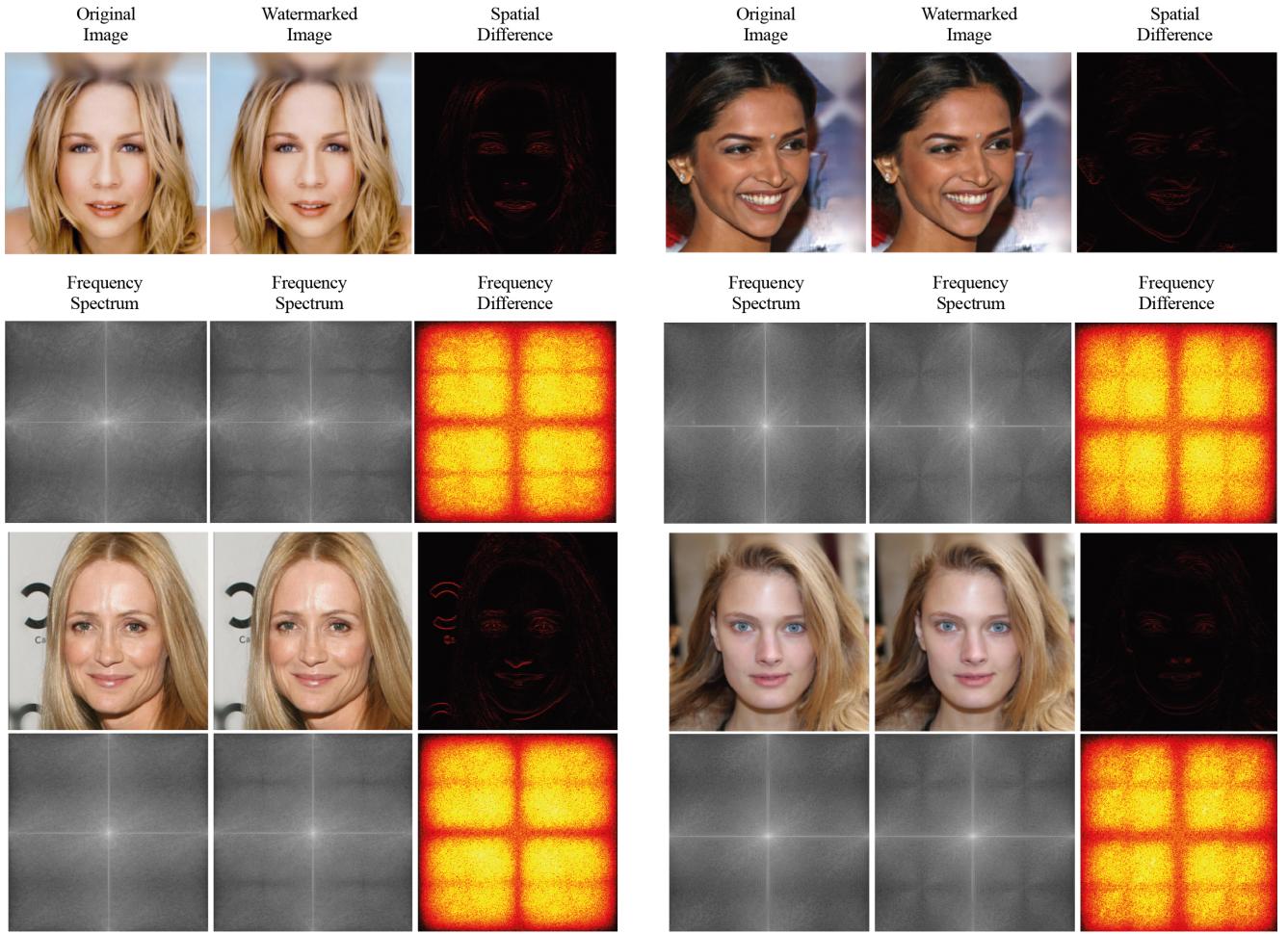


Figure 6. Visualization of the watermarking process using SpecGuard. The first row shows the original image, the watermarked image, and their spatial difference. The spatial difference highlights the minimal perceptual change between the original and watermarked images, ensuring imperceptibility. The second row presents the frequency spectrum of the original and watermarked images, along with their frequency difference, emphasizing the subtle embedding of the watermark in the high-frequency components. The comparison confirms that SpecGuard achieves invisible watermarking while maintaining robust frequency-domain characteristics for effective bit recovery.

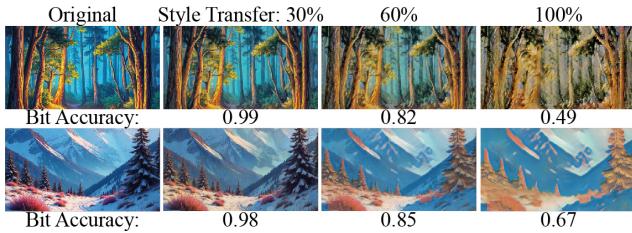


Figure 7. Effect of style transfer severity on bit recovery accuracy. As style intensity increases, bit accuracy decreases, showing the impact of major transformations.

692 A strength factor of $s(20)$ achieves optimal performance
 693 with a PSNR/SSIM of 42.89/0.99 and a BRA of 0.99. In-
 694 creasing s beyond 20 reduces PSNR and SSIM values, indi-
 695 cating diminished perceptual quality, while lower strength
 696 factors compromise robustness. Therefore, $s(20)$ effec-

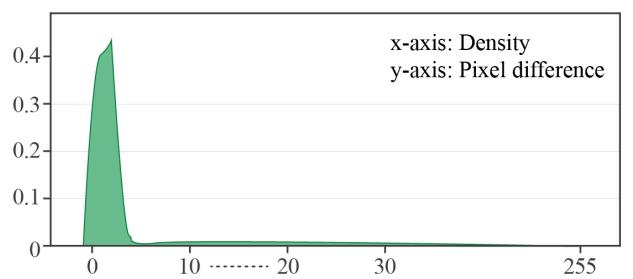


Figure 8. Pixel difference distribution between the original and watermarked images. The x-axis represents the pixel intensity difference, and the y-axis indicates the density. Most pixel differences remain close to zero, highlighting SpecGuard's minimal perceptual loss and superior imperceptibility of the embedded watermark.

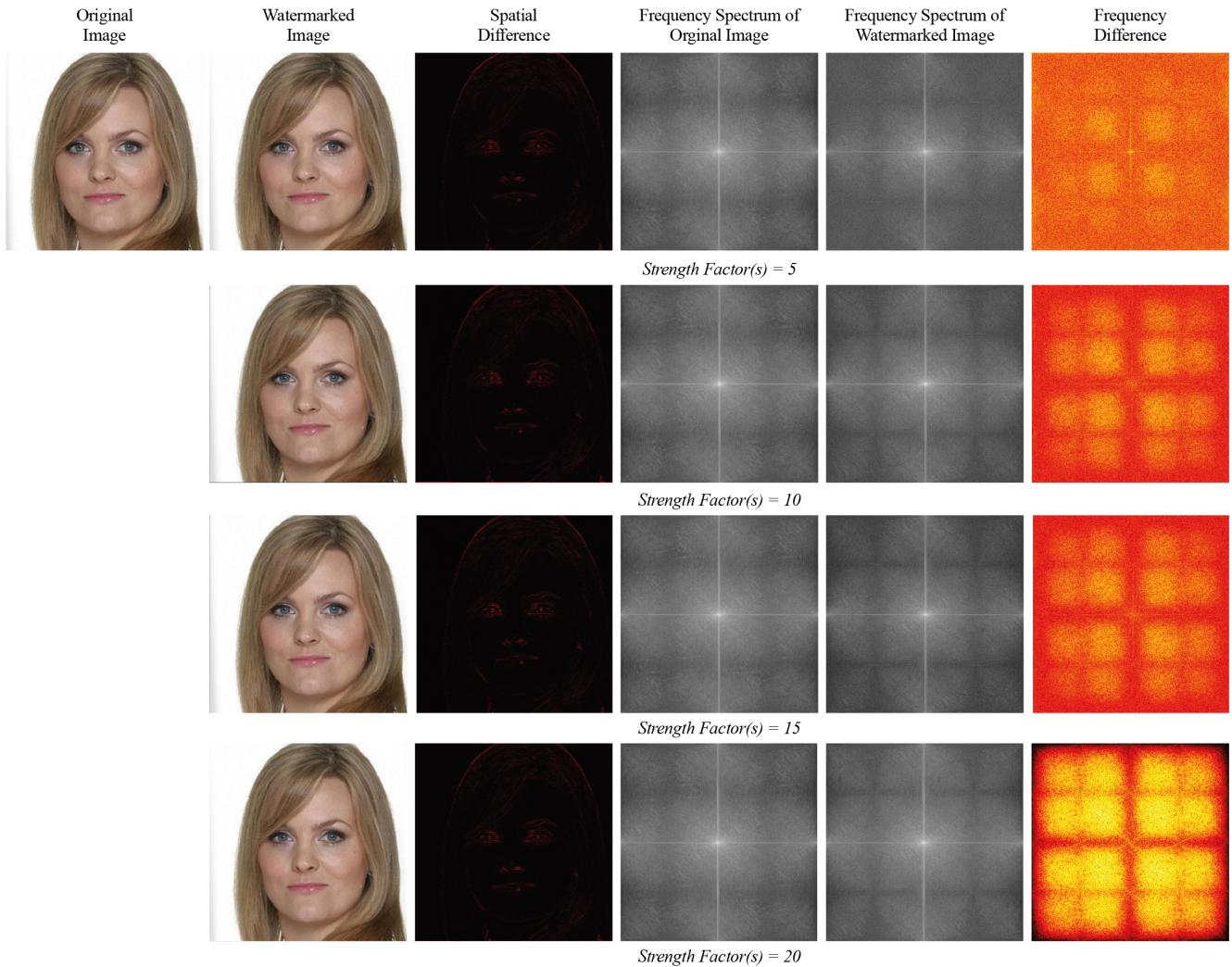


Figure 9. Visualization of the watermarking process using SpecGuard for different strength factors (s). The first row illustrates the original image, the watermarked image, and their spatial difference for $s = 5$, followed by the frequency spectra of the original and watermarked images and their frequency difference. The subsequent rows demonstrate the impact of increasing the strength factor ($s = 10, 15, 20$) on the frequency difference, highlighting the progressive embedding intensity. Higher strength factors increase the visibility in the frequency domain while maintaining imperceptibility in the spatial domain, ensuring robust watermarking without compromising image quality.

697 tively balances robustness and visual quality as also shown
698 in Fig. 6.

699 Figure 9 further demonstrates the effect of different
700 strength factors ($s = 5, 10, 15, 20$) on the watermark em-
701 bedding process. The first row showcases the original im-
702 age, the watermarked image, and their spatial difference,
703 highlighting the imperceptibility of the watermark in the
704 spatial domain. The subsequent rows compare the fre-
705 quency spectrum of the original and watermarked images,
706 as well as the frequency difference, illustrating how in-
707 creased strength factors enhance the visibility of the wa-
708 termark in the frequency domain while maintaining imper-
709 perceptibility in the spatial domain. Illustrate the robustness

710 and adaptability of the proposed SpecGuard model in em-
711 bedding and retaining watermark information under varying
712 conditions.

10. Description of Benchmarking Attacks

To comprehensively evaluate watermark robustness, we
714 benchmark performance against a diverse set of attacks, in-
715 cluding distortions, regeneration, and adversarial manipula-
716 tions. These attacks, derived from prior benchmarking ef-
717 forts [6], assess the stability of watermarks under real-world
718 transformations. The results are presented in Tab. 3 (main
719 paper) and the details of the attacks are in Tab. 9, compar-
720 ing multiple state-of-the-art (SOTA) methods such as Tree-
721

| Attack Name | Description | Parameters |
|-----------------------------|--|--|
| Distortion Attacks | | |
| Rotation | Rotates an image by a specified angle to test watermark robustness against geometric transformations. | Angle: 9° to 45° clockwise |
| Crop | Crops a portion of the image and resizes it back, simulating common editing. | Crop Ratio: 10% to 50% |
| Bright | Adjusts image brightness to test watermark stability under illumination changes. | Brightness Increase: 20% to 100% |
| Contrast | Modifies image contrast to simulate lighting variations. | Contrast Increase: 20% to 100% |
| Blur | Applies a low-pass filter to smooth the image, reducing high-frequency details. | Kernel Size: 4 to 20 pixels |
| Noise | Introduces random pixel fluctuations to simulate compression noise and low-quality rendering. | Std. Deviation: 0.02 to 0.1 |
| JPEG | Compresses the image using JPEG encoding, reducing quality and adding artifacts. | Quality Score: 90 to 10 |
| Geo | Combination of geometric distortion attacks, including rotation, crop, applied uniformly to assess cumulative effects. | Strength: Geo(x): Rotation: $9^\circ + x \times (45^\circ - 9^\circ)$, Crop: $10\% + x \times (50\% - 10\%)$ |
| Deg | Combination of degradation attacks, integrating blur, noise, and JPEG to simulate complex real-world distortions. | Strength: Deg(x): Blur: $4 + x \times (20 - 4)$, Noise: $0.02 + x \times (0.1 - 0.02)$, JPEG: $90 - x \times (90 - 10)$ |
| Regeneration Attacks | | |
| Regen-Diff | Passes an image through a diffusion model to reconstruct a similar but altered version. | Denoising Steps: 40 to 200 |
| Regen-DiffP | A prompted version of diffusion-based regeneration, leveraging text guidance to refine results. | Denoising Steps: 40 to 200 with Prompt |
| Regen-VAE | Uses a variational autoencoder to encode and decode an image, affecting watermark integrity. | Quality Level: 1 to 7 |
| Regen-KLVAE | Uses a KL-regularized autoencoder to compress and reconstruct an image, weakening watermark signals. | Bottleneck Sizes: 4, 8, 16, 32 |
| Rinse-2xDiff | Applies a two-stage diffusion regeneration, progressively altering the image over multiple steps. | Timesteps: 20 to 100 per diffusion |
| Rinse-4xDiff | Performs four cycles of diffusion-based image reconstruction, aggressively erasing watermark traces. | Timesteps: 10 to 50 per diffusion |
| Adversarial Attacks | | |
| AdvEmbG-KLVAE8 | Embeds adversarial perturbations using a grey-box VAE-based attack to reduce detection accuracy. | KL-VAE Encoding, $\epsilon = 2/255$ to $8/255$, PGD Iterations = 100, Step Size = $0.01 \times \epsilon$ |
| AdvEmbB-RN18 | Uses a pre-trained ResNet18 model to introduce adversarial noise and affect watermark recognition. | ℓ_∞ Perturbation: $2/255$ to $8/255$, PGD Iterations = 50, Step Size = $0.01 \times \epsilon$ |
| AdvEmbB-CLIP | Attacks the CLIP image encoder to introduce embedding shifts that disrupt watermark decoding. | ℓ_2 Perturbation Norm = 2.5, PGD Iterations = 50, Learning Rate = 0.001 |
| AdvEmbB-KLVAE16 | Uses an alternative KL-VAE model to introduce structured perturbations into the embedding process. | KL-VAE Embedding, Latent Size = 16, ℓ_∞ Perturbation = 4/255 |
| AdvEmbB-SdxlVAE | Attacks Stable Diffusion XL's VAE encoder to alter latent representations and remove watermarks. | Targeted VAE Perturbation, Diffusion Steps = 100, ℓ_2 Perturbation = 3.0 |
| AdvCls-UnWM&WM | Trains a surrogate detector on watermarked and non-watermarked images to bypass watermark detection. | Dataset Size = 3000 Images (1500 Per Class), ResNet-18, Learning Rate = 0.001, Batch Size = 128 |
| AdvCls-Real&WM | Trains an adversarial classifier using real and watermarked images to classify watermark presence. | Dataset Size = 15,000 Images (7500 Per Class), Adam Optimizer, Learning Rate = 0.0005, Batch Size = 128, Epochs = 10 |
| AdvCls-WM1&WM2 | Exploits watermark signal variations between different users to remove or alter hidden information. | Two Sets of Watermarked Images, Model = Vision Transformer (ViT), PGD Attack, Perturbation Strength = 6/255 |

Table 9. Overview of attack types, their mechanisms, and key parameters based on the prior study [6] that we also utilized in our study.

722 Ring [50], Stable Signature [34], and StegaStamp [44]. The
723 attacks are categorized as follows:

724 10.1. Distortion Attacks

725 These include standard image-processing transformations
726 that alter the spatial or color properties of images. We con-
727 sider rotation (9° to 45°) where images are rotated at vary-
728 ing degrees to test watermark stability. Resized cropping
729 (10% to 50%) removes portions of an image and resizes the
730 remaining content, mimicking common real-world editing.
731 Random erasing (5% to 25%) replaces regions with gray
732 pixels, simulating object removal. Brightness adjustments
733 (20% to 100%) and contrast modifications (20% to 100%)
734 simulate lighting variations. Gaussian blur (4 to 20 pixels)
735 applies low-pass filtering, while Gaussian noise (0.02 to 0.1
736 standard deviation) adds random pixel fluctuations, simu-
737 lating compression noise [6].

738 10.2. Regeneration Attacks

739 These attacks leverage generative models such as diffusion
740 and variational autoencoders (VAEs) to reconstruct images
741 while suppressing embedded watermarks. We evaluate sin-
742 ggle regeneration attacks including Regen-Diff (diffusion-
743 based reconstruction), Regen-DiffP (perceptually optimized
744 diffusion), Regen-VAE (autoencoder-based reconstruction),
745 and Regen-KLVAE (KL-regularized VAE reconstruction).
746 Additionally, multi-step regeneration attacks such as Rinse-
747 2xDiff and Rinse-4xDiff involve iterative diffusion pro-
748 cesses designed to further erase watermark traces [39, 57].

749 10.3. Adversarial Attacks

750 These attacks attempt to deceive watermark detectors
751 through embedding perturbations or surrogate model train-
752 ing. Grey-box embedding attacks (AdvEmbG-KLVAE8)
753 perturb watermarks while preserving image content. Black-
754 box embedding attacks (AdvEmbB-RN18, AdvEmbB-
755 CLIP, AdvEmbB-KLVAE16, AdvEmbB-SdxIVAE) intro-
756 duce noise during watermark embedding to decrease
757 detection confidence. Adversarial classifiers (AdvCls-
758 UnWM&WM, AdvCls-Real&WM, AdvCls-WM1&WM2)
759 use learned classifiers to distinguish watermarked images
760 and remove hidden signals [14, 35, 37, 39].

761 Overall, our evaluation framework ensures a rigorous as-
762 sessment of watermark robustness under various real-world
763 transformations and adversarial strategies.

764 References

- 765 [1] Edward H Adelson, Eero Simoncelli, and Rajesh Hingorani.
766 Orthogonal pyramid transforms for image coding. In *Vi-
767 sual Communications and image processing II*, pages 50–58.
768 SPIE, 1987. 2
- 769 [2] Mahdi Ahmadi, Alireza Norouzi, Nader Karimi, Shadrokh
770 Samavi, and Ali Emami. Redmark: Framework for resid-

- 771 ual diffusion watermarking based on deep networks. *Expert
772 Systems with Applications*, 146:113157, 2020. 2
- 773 [3] Aashutosh AV, Srijan Das, Abhijit Das, et al. Latent flow
774 diffusion for deepfake video generation. In *Proceedings of
775 the IEEE/CVF Conference on Computer Vision and Pattern
776 Recognition*, pages 3781–3790, 2024. 1
- 777 [4] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin
778 Hwang, and Nick Johnston. Variational image compression
779 with a scale hyperprior. *arXiv preprint arXiv:1802.01436*,
780 2018. 8
- 781 [5] Tali Dekel, Michael Rubinstein, Ce Liu, and William T Free-
782 man. On the effectiveness of visible watermarks. In *Pro-
783 ceedings of the IEEE Conference on Computer Vision and Pattern
784 Recognition*, pages 2146–2154, 2017. 2
- 785 [6] Mucong Ding, Tahseen Rabbani, Bang An, Aakriti Agrawal,
786 Yuancheng Xu, Chenghao Deng, Sicheng Zhu, Abdurisak
787 Mohamed, Yuxin Wen, Tom Goldstein, et al. Waves: Bench-
788 marking the robustness of image watermarks. In *ICLR 2024
789 Workshop on Reliable and Responsible Foundation Models*,
790 2024. 1, 6, 5, 7
- 791 [7] Hubert Etienne. The future of online trust (and why deepfake
792 is advancing it). *AI and Ethics*, 1(4):553–562, 2021. 1
- 793 [8] Han Fang, Yupeng Qiu, Kejiang Chen, Jiyi Zhang, Weim-
794 ing Zhang, and Ee-Chien Chang. Flow-based robust water-
795 marking with invertible noise layer for black-box distortions.
796 In *Proceedings of the AAAI conference on artificial intelli-
797 gence*, pages 5054–5061, 2023. 6
- 798 [9] Jianwei Fei, Zhihua Xia, Benedetta Tondi, and Mauro Barni.
799 Supervised gan watermarking for intellectual property pro-
800 tection. In *2022 IEEE International Workshop on Infor-
801 mation Forensics and Security (WIFS)*, pages 1–6. IEEE, 2022.
802 2
- 803 [10] S. A. Fulling. The local geometric asymptotics of contin-
804 uum eigenfunction expansions. ii. one-dimensional systems.
805 *SIAM Journal on Mathematical Analysis*, 14(4):605–623,
806 1983. 2
- 807 [11] Mahdieh Ghazvini, Elham Mohamadi Hachrood, and Mo-
808 jdeh Mirzadi. An improved image watermarking method in
809 frequency domain. *Journal of Applied Security Research*, 12
810 (2):260–275, 2017. 2
- 811 [12] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep
812 sparse rectifier neural networks. In *Proceedings of the four-
813 teenth international conference on artificial intelligence and
814 statistics*, pages 315–323. JMLR Workshop and Conference
815 Proceedings, 2011. 3
- 816 [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing
817 Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and
818 Yoshua Bengio. Generative adversarial nets. *Advances in
819 neural information processing systems*, 27, 2014. 1
- 820 [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
821 Deep residual learning for image recognition. In *Pro-
822 ceedings of the IEEE conference on computer vision and pattern
823 recognition*, pages 770–778, 2016. 7
- 824 [15] K Jayashre and M Amsaprabhaa. Safeguarding media in-
825 tegrity: A hybrid optimized deep feature fusion based deep-
826 fake detection in videos. *Computers & Security*, 142:103860,
827 2024. 1

- 828 [16] Zhaoyang Jia, Han Fang, and Weiming Zhang. Mbrs: En- 829 hancing robustness of dnn-based watermarking by mini- 830 batch of real and simulated jpeg compression. In *Proceed- 831 ings of the 29th ACM international conference on multime- 832 dia*, pages 41–49, 2021. 6
- 833 [17] Hoon Kang and Joonsoo Ha. Projection spectral analysis. 834 *International Journal of Control, Automation and Systems*, 835 13(6):1530–1537, 2015. 1, 2, 4
- 836 [18] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 837 Progressive growing of gans for improved quality, stability, 838 and variation. In *International Conference on Learning Rep- 839 resentations (ICLR)*, 2018. 6
- 840 [19] SS Kelkar, LL Grigsby, and J Langsner. An extension of 841 parseval’s theorem and its use in calculating transient energy 842 in the frequency domain. *IEEE Transactions on Industrial 843 Electronics*, (1):42–45, 1983. 2, 5, 1
- 844 [20] Varsha Kishore, Xiangyu Chen, Yan Wang, Boyi Li, and 845 Kilian Q Weinberger. Fixed neural network steganography: 846 Train the images, not the network. In *International Confer- 847 ence on Learning Representations*, 2021. 2
- 848 [21] Jae-Eun Lee, Young-Ho Seo, and Dong-Wook Kim. Con- 849 volutional neural network-based digital image watermarking 850 adaptive to the resolution of image and watermark. *Applied 851 Sciences*, 10(19):6854, 2020. 2
- 852 [22] Yicheng Leng, Chaowei Fang, Gen Li, Yixiang Fang, and 853 Guanbin Li. Removing interference and recovering content 854 imaginatively for visible watermark removal. In *Proceed- 855 ings of the AAAI Conference on Artificial Intelligence*, pages 856 2983–2990, 2024. 2
- 857 [23] Xiao Li, Liquan Chen, Ju Jia, Zhongyuan Qin, and Zhangjie 858 Fu. A lightweight image forgery prevention scheme for iot 859 using gan-based steganography. *IEEE Transactions on Indus- 860 trial Informatics*, 2024. 2
- 861 [24] Dongdong Lin, Benedetta Tondi, Bin Li, and Mauro Barni. 862 A cyclegan watermarking method for ownership verification. 863 *IEEE Transactions on Dependable and Secure Computing*, 864 2024. 2
- 865 [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, 866 Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence 867 Zitnick. Microsoft coco: Common objects in context. In 868 *Computer Vision–ECCV 2014: 13th European Conference, 869 Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 5, 6
- 871 [26] Yang Liu, Zhen Zhu, and Xiang Bai. Wdnet: Watermark- 872 decomposition network for visible watermark removal. In 873 *Proceedings of the IEEE/CVF winter conference on applica- 874 tions of computer vision*, pages 3685–3693, 2021. 2
- 875 [27] Xiyang Luo, Ruohan Zhan, Huiwen Chang, Feng Yang, and 876 Peyman Milanfar. Distortion agnostic deep watermarking. 877 In *Proceedings of the IEEE/CVF conference on computer vi- 878 sion and pattern recognition*, pages 13548–13557, 2020. 2
- 879 [28] Thanh Thi Nguyen, Quoc Viet Hung Nguyen, Dung Tien 880 Nguyen, Duc Thanh Nguyen, Thien Huynh-The, Saeid 881 Nahavandi, Thanh Tam Nguyen, Quoc-Viet Pham, and 882 Cuong M Nguyen. Deep learning for deepfakes creation and 883 detection: A survey. *Computer Vision and Image Under- 884 standing*, 223:103525, 2022. 1
- 29 [29] Guangyu Nie, Changhoon Kim, Yezhou Yang, and Yi Ren. 885 Attributing image generative models using latent finger- 886 prints. In *International Conference on Machine Learning*, 887 pages 26150–26165. PMLR, 2023. 2
- 30 [30] Li Niu, Xing Zhao, Bo Zhang, and Liqing Zhang. Fine- 888 grained visible watermark removal. In *Proceedings of the 889 IEEE/CVF International Conference on Computer Vision*, 890 pages 12770–12779, 2023. 2
- 31 [31] Konstantin A Pantserov. The malicious use of ai-based deep- 891 fake technology as the new threat to psychological security 892 and political stability. *Cyber defence in the age of AI, smart 893 societies and augmented humanity*, pages 37–55, 2020. 1
- 32 [32] Ram Shankar Pathak. *The wavelet transform*. Springer Sci- 894 ence & Business Media, 2009. 1, 2, 4
- 33 [33] Fernandez Pierre, Alexandre Sablayrolles, Teddy Furon, 895 Hervé Jégou, and Matthijs Douze. Watermarking images in 896 self-supervised latent spaces. In *ICASSP 2022-2022 IEEE 897 International Conference on Acoustics, Speech and Signal 898 Processing (ICASSP)*, pages 3054–3058. IEEE, 2022. 2
- 34 [34] Fernandez Pierre, Guillaume Couairon, Hervé Jégou, 899 Matthijs Douze, and Teddy Furon. The stable signature: 900 Rooting watermarks in latent diffusion models. In *Proceed- 901 ings of the IEEE/CVF International Conference on Com- 902 puter Vision*, pages 22466–22477, 2023. 1, 6, 7
- 35 [35] Dustin Podell, Zion English, Kyle Lacey, Andreas 902 Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and 903 Robin Rombach. Sdxl: Improving latent diffusion mod- 904 els for high-resolution image synthesis. *arXiv preprint 905 arXiv:2307.01952*, 2023. 7
- 36 [36] Tong Qiao, Yuyan Ma, Ning Zheng, Hanzhou Wu, Yanli 906 Chen, Ming Xu, and Xiangyang Luo. A novel model water- 907 marking for protecting generative adversarial network. *Com- 908 puters & Security*, 127:103102, 2023. 2
- 37 [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 909 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, 910 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning 911 transferable visual models from natural language supervi- 912 sion. In *International conference on machine learning*, pages 913 8748–8763. PMLR, 2021. 7
- 38 [38] David E Rumelhart, Geoffrey E Hinton, and Ronald J 914 Williams. Learning internal representations by error prop- 915 agation, parallel distributed processing, explorations in the 916 microstructure of cognition, ed. de rumelhart and j. mcclel- 917 land. vol. 1. 1986. *Biometrika*, 71(599–607):6, 1986. 3
- 39 [39] Mehrdad Saberi, Vinu Sankar Sadasivan, Keivan Rezaei, 918 Aounon Kumar, Atoosa Chegini, Wenxiao Wang, and Soheil 919 Feizi. Robustness of ai-image detectors: Fundamental lim- 920 its and practical attacks. *arXiv preprint arXiv:2310.00076*, 921 2023. 7
- 40 [40] Sunpreet Sharma, Ju Jia Zou, Gu Fang, Pancham Shukla, and 922 Weidong Cai. A review of image watermarking for identity 923 protection and verification. *Multimedia Tools and Applica- 924 tions*, 83(11):31829–31891, 2024. 2
- 41 [41] Qingtang Su, Huanying Wang, DeCheng Liu, Zihan Yuan, 925 and Xueteng Zhang. A combined domain watermarking al- 926 gorithm of color image. *Multimedia Tools and Applications*, 927 79(39):30023–30043, 2020. 2

- 942 [42] Qingtang Su, Xuetong Zhang, and Huanying Wang. A blind 999
943 color image watermarking algorithm combined spatial domain 1000
944 and svd. *International Journal of Intelligent Systems*, 33(8):4747–4771, 2022. 2
945
- 946 [43] Ruizhou Sun, Yukun Su, and Qingyao Wu. Denet: disentangled 1001
947 embedding network for visible watermark removal. In *Proceedings 1002
948 of the AAAI Conference on Artificial Intelligence*, pages 2411–2419, 2023. 2
949
- 950 [44] Matthew Tancik, Ben Mildenhall, and Ren Ng. Stegastamp: 1003
951 Invisible hyperlinks in physical photographs. In *Proceedings 1004
952 of the IEEE/CVF conference on computer vision and pattern 1005
953 recognition*, pages 2117–2126, 2020. 1, 6, 7
- 954 [45] D Vaishnavi and TS Subashini. Robust and invisible image 1006
955 watermarking in rgb color space using svd. *Procedia Computer 1007
956 Science*, 46:1770–1777, 2015. 2
- 957 [46] Guanjie Wang, Zehua Ma, Chang Liu, Xi Yang, Han Fang, 1008
958 Weiming Zhang, and Nenghai Yu. Must: Robust image 1009
959 watermarking for multi-source tracing. In *Proceedings of 1010
960 the AAAI Conference on Artificial Intelligence*, pages 5364– 1011
961 5371, 2024. 6
- 962 [47] Huanying Wang and Qingtang Su. A color image water- 1012
963 marking method combined qr decomposition and spatial 1013
964 domain. *Multimedia Tools and Applications*, 81(26):37895–
965 37916, 2022. 2
- 966 [48] Zilan Wang, Junfeng Guo, Jiacheng Zhu, Yiming Li, Heng 1014
967 Huang, Muhan Chen, and Zhengzhong Tu. Sleepermark: 1015
968 Towards robust watermark against fine-tuning text-to-image 1016
969 diffusion models. *arXiv preprint arXiv:2412.04852*, 2024. 6
- 970 [49] Zijie J Wang, Evan Montoya, David Munechika, Haoyang 1017
971 Yang, Benjamin Hoover, and Duen Horng Chau. Diffusiondb: 1018
972 A large-scale prompt gallery dataset for text-to- 1019
973 image generative models. *arXiv preprint arXiv:2210.14896*, 1020
974 2022. 5
- 975 [50] Yuxin Wen, John Kirchenbauer, Jonas Geiping, and Tom 1021
976 Goldstein. Tree-rings watermarks: Invisible fingerprints for 1022
977 diffusion images. In *Advances in Neural Information Pro- 1023
978 cessing Systems*, pages 58047–58063. Curran Associates,
979 Inc., 2023. 1, 6, 7
- 980 [51] Yuxin Wen, John Kirchenbauer, Jonas Geiping, and Tom 1024
981 Goldstein. Tree-rings watermarks: Invisible fingerprints for 1025
982 diffusion images. *Advances in Neural Information Process- 1026
983 ing Systems*, 36, 2024. 1
- 984 [52] Bing Xu. Empirical evaluation of rectified activations in con- 1027
985 volutional network. *arXiv preprint arXiv:1505.00853*, 2015. 3
- 986
- 987 [53] Zijin Yang, Kai Zeng, Kejiang Chen, Han Fang, Weim- 1028
988 ing Zhang, and Nenghai Yu. Gaussian shading: Prov- 1029
989 able performance-lossless image watermarking for diffusion 1030
990 models. In *Proceedings of the IEEE/CVF Conference on 1031
991 Computer Vision and Pattern Recognition*, pages 12162– 1032
992 12171, 2024. 1, 6
- 993 [54] Zihan Yuan, Qingtang Su, Decheng Liu, and Xuetong Zhang. 1033
994 A blind image watermarking scheme combining spatial do- 1034
995 main and frequency domain. *The visual computer*, 37:1867– 1035
996 1881, 2021. 2
- 997 [55] Chaoning Zhang, Philipp Benz, Adil Karjauv, Geng Sun, and 1036
998 In So Kweon. Udh: Universal deep hiding for steganography,
999