# ENHANCED SCREEN SHOOTING RESILIENT DOCUMENT WATERMARKING

*Heng Wang*    *Hongxia Wang**    *Xinyi Huang*    *Zhenhao Shi*

School of Cyber Science and Engineering, Sichuan University, Chengdu, China;
Key Laboratory of Data Protection and Intelligent Management (Sichuan University),
Ministry of Education, China.

## ABSTRACT

The widespread adoption of smartphones has introduced new challenges to document copyright protection, prompting the emergence of Screen-Shooting Resilient Document Watermarking (SSRDW) technology. In recent years, underpainting-based SSRDW techniques have proven to be highly effective. However, after careful study, we find that existing methods fail to simultaneously meet four essential criteria for SSRDW: high imperceptibility, strong robustness, adaptability to text processing, and high efficiency. In this paper, we introduce an enhanced underpainting-based SSRDW approach capable of satisfying all four requirements. Our approach enhances imperceptibility by employing underpainting embedding methods independent of text content. Additionally, we introduce a fast resynchronization mechanism to improve time efficiency. Furthermore, we propose an enhanced watermark extraction method that enhances robustness and enables watermark retrieval even in scenarios involving text processing. Extensive experimental validation underscores the superior performance of our enhanced SSRDW method.

*Index Terms*— Document watermarking, Imperceptibility, Robustness, Adaptability, Efficiency

## 1. INTRODUCTION

Documents are widely used for information storage in various fields. However, with the widespread adoption of smartphones, photography has emerged as a straightforward and efficient mode of information transmission. The emergence of this trend has prompted an urgent demand for document anti-capture measures, leading to the development of screen-shooting resilient document watermarking (SSRDW) technology. However, research in this field presents substantial challenges. Documents lack the auditory or visual characteristics found in media such as audio [1, 2] or images [3], and the current state of natural language processing research has limitations in comprehending text content.

Current document watermarking technology are typically categorized into two main types: text layer-based and graphical layer-based approaches. Text layer-based approaches encompass text structure [4, 5], font structure [6, 7, 8], and semantics [9, 10] techniques, but they often entail modifications to the text content, thereby imposing certain constraints. In graphical layer-based design, the focus is primarily on crafting the underpainting of the document.

This approach, not constrained by text content, has demonstrated its utmost effectiveness in document watermarking.

In this paper, we primarily focus on the critical aspects of the SSRDW process, which include high imperceptibility, strong robustness, adaptability to text processing, and high efficiency. Imperceptibility entails the imperative to preserve the visual quality of watermarked documents to render them inconspicuous to the human eye. Robustness signifies the document's capacity to extract watermark information even following image capture under various conditions such as different distances and angles. Adaptability ensures that the watermark can still be extracted despite text processing adjustments. Lastly, efficiency refers to the need for prompt responses in the embedding and extraction processes.

However, after careful investigation, it was determined that none of the underpainting-based methods could simultaneously fulfill all four aforementioned criteria. Masahiko *et al.* [11] proposed using distinct dot arrays to represent different watermark signals, Pramila *et al.* [12] employed periodic templates for watermark encoding. Nevertheless, these methods lacked synchronization mechanisms, thereby compromising robustness during partial capture. Gugelmann *et al.* [13] introduced a watermark embedding method by modify underpainting brightness, yet this approach suffered from reduced visual quality, failing to ensure imperceptibility. Fang *et al.* [14] presented a multiple watermark embedding scheme, but it exhibited lower time efficiency during resynchronization and required a strong association between underpainting and text, thus limiting its adaptibility. To address these limitations, we present a novel enhanced underpainting-based SSRDW scheme. To the best of our knowledge, this is the first underpainting-based solution capable of simultaneously fulfilling all four requirements. Our contributions are as follows:

- We devise an independent underpainting embedding method, enabling underpainting to be embedded without constraints imposed by text formatting, thus further enhancing imperceptibility and adaptability.

- We address non-uniform distortion in screen capture with a partially overlapped sub-block histogram-equalization pre-processing approach and introduce a fast synchronization mechanism to enhance time efficiency.

- We propose an enhanced watermark extraction approach that not only improves robustness but also ensures watermark retrieval in scenarios involving text processing, greatly enhancing adaptability.

## 2. PROPOSED METHOD

The overall framework, as shown in Fig. 1, is primarily divided into two parts: the watermark embedding and the watermark extraction. Detailed descriptions of these components will follow.
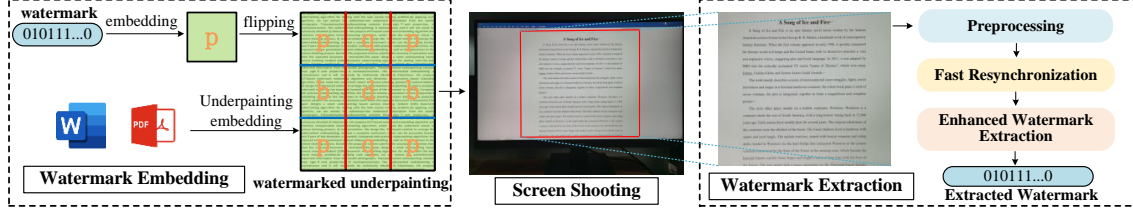
**Fig. 1**. Overall framework of the system, which is divided into two parts: the watermark embedding and watermark extraction.

## 2.1. Watermark Embedding

The watermark embedding primarily consists of two steps: underpainting generation and underpainting flipping.

**1) Underpainting Generation:** The watermark bits are embedded based on the size relationship of a pair of DCT coefficients. Specifically, for a $30 \times 30$ pixels block, the discrete cosine transform (DCT) is initially applied, resulting in a $30 \times 30$ DCT matrix $\mathbf{C}$, the mid-frequency coefficients $c_1 = \mathbf{C}(4,5)$ and $c_2 = \mathbf{C}(5,4)$ are selected. The embedding process can be expressed as follows:

$$\begin{cases} c_1 = r, c_2 = -r, & if \ w = 0 \\ c_1 = -r, c_2 = r, & if \ w = 1 \end{cases} \quad (1)$$

where $w \in \{0, 1\}$, representing the watermark bit, and $r \in \mathbb{N}^*$, denotes the embedding strength. Thus, for a watermark information of $a \times b$ bits, a underpainting of $a \times b \times (30 \times 30)$ pixels is required. This results a complete watermark unit, indicated as $\boldsymbol{p}$ in Fig. 1.

**2) Underpainting Flipping:** After obtaining $\boldsymbol{p}$, Fang *et al.* [14] performed underpainting scaling based on the text, which often reduced adaptability due to the strong association between underpainting and text. Here, we direct perform row and column symmetry construction to obtain the watermarked underpainting shown in Fig. 1. This strategic adjustment, combined with the implementation of an enhanced watermark extraction approach, ensures robustness during the extraction process.

This watermarked underpainting contains several complete watermark units, with the red and blue lines representing the boundaries of a complete unit, which are also the targets to be detected in the resynchronization phase.

## 2.2. Watermark Extraction

For captured image $\boldsymbol{I}$, we begin by delineating a quadrilateral region based on the document distribution. This region is then subjected to perspective transformation to attain a $2048 \times 2048$ resolution image denoted as $\boldsymbol{I'}$. Subsequently, we perform preprocessing, fast resynchronization, and enhanced watermark extraction.

**1) Preprocessing**: Screen-shooting process is often accompanied by a series of distortions that subsequently affect the subsequent resynchronization. Fang *et al.* [14] utilized a non-overlapping sub-block histogram equalization approach for preprocessing, but due to the non-uniform nature of distortion, this approach often lacks treatment for transition between sub-blocks. Partially overlapped sub-block histogram-equalization (POSHE) [15], as an image enhancement method, significantly reduces variations in the equalization functions between adjacent sub-blocks. Thus, we utilize POSHE for preprocessing with four steps: 1) Initialize the sub-block size and step size, 2) Perform histogram equalization on the current sub-block, move to the next block according to the step size, 3) Repeat step 2) until the entire $\boldsymbol{I'}$ is covered, and 4) Divide each pixel in the resulting image by the frequency of sub-block histogram equalization. This yields the preprocessed image $\boldsymbol{I}_p$.

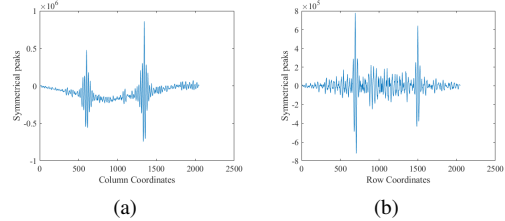We opted for a 16-pixel sub-block size with an 8-pixel step size.



**Fig. 2**. Symmetric waveform detection results: (a) Column symmetric waveform. (b) Row symmetric waveform.

This simultaneous processing by multiple sub-blocks enhances treatment of transition regions and fine underpainting details, significantly boosting resynchronization accuracy.

**2) Fast Resynchronization**: As motioned before, we first need to locate complete watermark regions in $\boldsymbol{I}_p$. Unlike the pixel-wise calculation approach in [14], we have devised a fast resynchronization method based on the frequency-domain characteristics of the symmetrical underpainting. This method consists of two stages: column detection and row detection. Taking column detection as an example, we initially partition $\boldsymbol{I}_p$ into $32 \times 32$ non-overlapping image blocks. Subsequently, we perform the following normalization operation on each block:

$$C(B) = \frac{B - E(B)}{\sqrt{D(B)}} \quad (2)$$

where $B$ represents each block, $E(B)$ denotes the mean of $B$, and $D(B)$ represents the variance of $B$. Consequently, we obtain a matrix $\boldsymbol{I}_{nor}$. We then generate $\boldsymbol{I}_{plr}$ by horizontally flipping $\boldsymbol{I}_{nor}$. The column symmetry matrix, denoted as $\boldsymbol{T}_{plr}$, is defined as the cross-correlation between $\boldsymbol{I}_{nor}$ and $\boldsymbol{I}_{plr}$ as follows:

$$\boldsymbol{T}_{plr} = \mathcal{D}(IFFT[FFT(\hat{\boldsymbol{I}_{nor}})FFT(\hat{\boldsymbol{I}_{plr}})^*]) \quad (3)$$

where $\hat{\boldsymbol{I}_{nor}}$ and $\hat{\boldsymbol{I}_{plr}}$ are obtained by zero-padding $\boldsymbol{I}_{nor}$ and $\boldsymbol{I}_{plr}$ to twice their original dimensions, respectively. *FFT* represent the fast Fourier transform, *IFFT* is the inverse transform, and * denotes complex conjugation. $\mathcal{D}(\cdot)$ is a downsampling function that reduces the input size by half. The final column symmetry $\boldsymbol{S}_{col} = \boldsymbol{T}_{plr}[2048, :]$ is given by Fig. 2(a). The top two peaks in this waveform represent the column coordinates of a complete watermark region, indicating the positions of the red lines on the watermarked underpainting shown in Fig. 1.

For row localization, the process is often affected by text interference. Here, we first extract the text regions from $\boldsymbol{I'}$ through binarization, obtaining $\boldsymbol{I}_{text}$. Subsequently, we perform vertical flipping on both $\boldsymbol{I}_{nor}$ and $\boldsymbol{I}_{text}$ to obtain $\boldsymbol{I}_{pud}$ and $\boldsymbol{I}_{text,pud}$. We then apply a procedure similar to column detection to these two matrices. Consequently, we derive $\boldsymbol{T}_{pud}$ and $\boldsymbol{T}_{text,pud}$. The final row symmetry $\boldsymbol{S}_{row} = \boldsymbol{T}_{pud}[:, 2048] - \boldsymbol{T}_{text,pud}[:, 2048]$, depicted in Fig. 2(b). The highest two peak in this waveform represents the row coordinates of a complete watermark region, corresponding to the positions of the blue lines on the watermarked underpainting shown in Fig. 1.

**Enhanced Watermark Extraction:** After obtaining a complete watermark region $R$, it is important to address the irreversible distortion caused by text overlay on the underpainting. A text region compensation method was introduced in [14], where the core idea is to compensate for the text region using inter-line spacing due to the symmetry exhibited by the watermark blocks generated by Eq. (1). However, this approach exhibits reduced accuracy in cases of substantial text coverage. Therefore, we propose an enhanced watermark extraction method. First, $R$ is resized to $a \times b \times (64 \times 64)$ pixels and segmented into $a \times b$ non-overlapping blocks. For each block $D$, the following operation is performed:

$$
\begin{cases}
D(x,y) = \begin{cases} D(64-x, 64-y), & if\, D'(x,y) = 0 \\ & \&\, D'(64-x, 64-y) = 1 \\ E(D), & otherwise \end{cases} \\
D_{en}(x,y) = \dfrac{1}{3} \sum_{i=0}^{2} \text{Stretch}(D(x,y), 2^i) \\
w' = \begin{cases} 0, & if \quad c_1 \geq c_2 \\ 1, & otherwise \end{cases}
\end{cases}
\tag{4}
$$

where $D'$ is obtained by binarization of $D$, $E(D)$ denotes the mean of $D$. $\text{Stretch}(D(x,y), 2^i)$ signifies the operation of partitioning $D$ into $2^i$ blocks and applying contrast stretching to each block. For enhanced block $D_{en}$, we then resize it into $30 \times 30$ pixels and extract watermark bit $w'$ by comparing its DCT coefficients $c_1$ and $c_2$. This method enables multi-scale analysis of text coverage, enhancing the underpainting texture. Fig. 3 further illustrates its effectiveness, where Fig. 3(a) show the compensated and enhanced watermark blocks, and Fig. 3(b) and Fig. 3(c) display the corresponding DCT coefficient matrices, with brighter regions indicating larger magnitude. It can be observed that the enhanced watermark block exhibits stronger discriminability at coordinates (4,5) and (5,4), thus facilitating a more favorable decoding outcome.
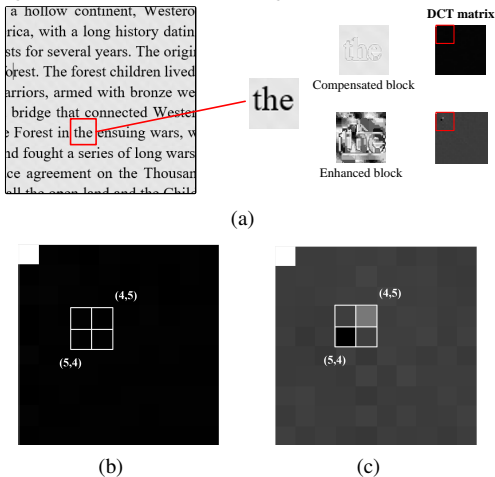


(a)



(b)                                    (c)

**Fig. 3**. Illustration of enhanced watermark block. (a) Compensated and enhanced watermark blocks. (b) DCT matrix for compensated watermark block. (c) DCT matrix for enhanced watermark block.

## 3. EXPERIMENTAL RESULTS

### 3.1. Experimental Setup

**Comparison algorithm.** We compared our method with three state-of-the-art watermarking algorithms (Pramila *et al.* [12], Gugelmann *et al.* [13], Fang *et al.* [14]). Among these, [12] is designed for

print-camera process, while [13] and [14] are designed for screen-shooting process.

**Test set**. The test dataset consists of 10 distinct documents, with content randomly excerpted from the Chinese book "The Three-Body Problem" [16] and the English book "A Song of Ice and Fire" [17], forming the *Dataset*.

**Metrics**. We utilize Peak Signal-to-Noise Ratio (PSNR) and Structural SIMilarity (SSIM) to evaluate imperceptibility, Robustness is evaluated through average erroneous bits (AEB), calculated as the average number of erroneous bits across all documents in the *Dataset*. Higher PSNR and SSIM values, coupled with lower AEB, indicate superior results.

**Parameters**. The error correction code we employ is Bose Chaudhuri Hocquenghem (BCH) code (64,36), capable of correcting up to 5 bits of errors. Within this coding scheme, 36 bits serve as the payload, accommodating a whopping $2^{36}$=68719476736 devices. This capacity is more than sufficient for a company's needs. And the embedding strength $r = 50$, corresponding PSNR=38.58dB.

**Environments**. The experiments were conducted using a "HUAWEI MATE 30" capturing device and a "PHL 275S9" display monitor. All experiments were carried out on a computer system featuring a 3.79 GHz Intel Core i7 CPU with 32GB of RAM, operating on a 64-bit Windows 10 platform.
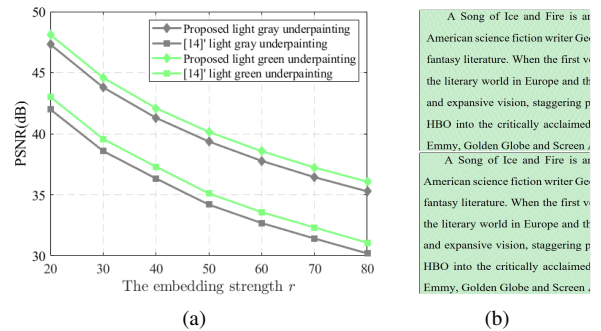


(a)                                    (b)

**Fig. 4**. Visual quality comparison between [14] and the proposed method. (a) PSNR of different embedding strengths. (b) *Top Row:* The embedded document with [14]. *Bottom Row:* The embedded document with proposed method.

**Table 1**. Visual quality comparison of different methods.

| Method | Pramila *et al.* [12] | Gugelmann *et al.* [13] | Fang *et al.* [14] | Proposed |
|---|---|---|---|---|
| document | world in Europe a ive vision, staggeri he critically accla den Globe and Scr | world in Europe a ve vision, staggeri he critically accla den Globe and Scr | world in Europe a ive vision, staggeri he critically accla den Globe and Scr | world in Europe a ve vision, staggeri he critically accla den Globe and Scr |
| PSNR(dB) | 38.33 | 37.89 | 38.12 | 38.58 |
| SSIM | 0.8647 | 0.8541 | 0.8952 | 0.9097 |

### 3.2. Imperceptibility

To evaluate imperceptibility, we initially compare our method with a representative screen-shooting document method [14]. We select a light gray color (RGB-[204, 204, 204]) and a light green color (RGB-[201, 230, 204]) as underpainting' colors. To ensure a fair comparison, both our proposed method and [14] were tested under the same parameters. As shown in Fig. 4(a), the proposed method exhibit higher PSNR value on both underpainting's colors compared

to [14]. We attribute this phenomenon to the fact that [14] scaled watermark blocks based on the text width, leading to more pronounced "block artifacts". In contrast, our approach, which avoids text region adaptation, results in improved PSNR. Fig. 4(b) illustrates the subjective visual comparison at the same embedding strength of $r = 50$.

Furthermore, as depicted in Table 1, we also compared our method with two other methods. To ensure a fair comparison of the subjective visual quality and robustness of different methods, we harmonized the embedding strength of all schemes to achieve comparable PSNR and SSIM values.

### 3.3. Robustness

We evaluate robustness in terms of shooting distance and angle. We conduct tests using a light gray (RGB-[204, 204, 204]) underpainting for all methods while aligning watermark bits.

**Table 2**. Comparison of the AEB at different shooting distances.

| Distance (cm) | Pramila et al. [12] | Gugelmann et al. [13] | Fang et al. [14] | Proposed |
|---|---|---|---|---|
| 35 | 5.3 | 8.2 | 5.0 | **3.0** |
| 45 | 5.7 | 10.2 | 4.5 | **2.2** |
| 55 | 6.8 | 10.3 | 3.8 | **2.1** |
| 65 | 8.9 | 12.1 | 4.6 | **2.0** |
| Average | 6.7 | 10.2 | 4.5 | **2.3** |

**Different Distances:** At distances ranging from 35 cm to 65 cm with 10 cm intervals, our approach outperforms [12] and [13], as demonstrated in Table 2. Furthermore, in comparison to [14], our method achieves a significant reduction in AEB, averaging only 2.3 bits. Across various capture distances, the proposed approach consistently keeps the number of erroneous bits well within the correctable range of BCH coding.

**Table 3**. Captured document images under different horizontal angles, along with perspective transformation images.

| horizontal angles(°) | Left 45 | Left 15 | Right 15 | Right 45 |
|---|---|---|---|---|
| Captured | | | | |
| Corrected | | | | |

**Table 4**. Comparison of the AEB at different horizontal angles.

| horizontal angles(°) | Pramila et al. [12] | Gugelmann et al. [13] | Fang et al. [14] | Proposed |
|---|---|---|---|---|
| Left 60 | 9.2 | 8.6 | **4.0** | **4.0** |
| Left 45 | 8.8 | 10.7 | 4.7 | **2.0** |
| Left 30 | 8.2 | 11.6 | 4.8 | **1.0** |
| Left 15 | 7.3 | 10.8 | 4.2 | **1.0** |
| Right 15 | 8.0 | 8.5 | 4.9 | **1.6** |
| Right 30 | 9.1 | 7.6 | 4.7 | **1.2** |
| Right 45 | 9.6 | 8.3 | 4.7 | **2.2** |
| Right 60 | 11.6 | 6.9 | 5.0 | **4.2** |
| Average | 9.0 | 9.1 | 4.6 | **2.2** |

**Different Horizontal Angles:** We conduct tests with the capturing device fixed at a 60 cm distance from the screen, spanning horizontal angles from left 60° to right 60° in 15° increments. Prior to watermark extraction, we applied perspective transformations based on the document's distribution (Table 3). Table 4 shows the AEB at various horizontal angles, with our method consistently achieving

the lowest AEB, averaging 4.6 to 2.2 bits less compared to [14]. It's worth noting that we consider vertical captures similar to horizontal ones, and due to space constraints, we did not discuss vertical captures.

### 3.4. Adaptability

To validate adaptability to text processing, we conduct tests by editing the text with varying line spacing (in pt). Table 5 illustrates the degree of underpainting coverage at different line spacings. We utilize the Times New Roman font at 12-point size and maintain a fixed distance of 60 cm for capturing. As shown in Table 6, our approach consistently achieved the lowest AEB across different line spacings. This indicates the effectiveness of our proposed method in enhancing adaptability.

**Table 5**. Coverage of underpainting at different line spacings.

| Line spacing (pt) | 10 | 12 | 16 | 20 |
|---|---|---|---|---|
| Degree of underpainting coverage | continent about th years ago. Each se this continent were nature and used m skills, landed in W continent (destroy Dronish Islands ar the Forest. The w under which the F | continent about th years ago. Each se this continent were nature and used m skills, landed in W continent (destroy Dronish Islands ar the Forest. The w | continent about th years ago. Each se this continent were nature and used m skills, landed in W continent (destroy | continent about th years ago. Each se this continent were nature and used m skills, landed in W |

**Table 6**. Comparison of the AEB at different line spacing.

| Line Spacing (pt) | Pramila et al. [12] | Gugelmann et al. [13] | Fang et al. [14] | Proposed |
|---|---|---|---|---|
| 10 | 28.6 | 20.6 | 15.3 | **3.6** |
| 12 | 24.1 | 18.1 | 11.3 | **3.2** |
| 14 | 18.9 | 14.6 | 8.1 | **2.6** |
| 16 | 15.5 | 10.1 | 6.5 | **2.5** |
| 18 | 12.3 | 9.8 | 4.3 | **2.1** |
| 20 | 9.1 | 9.6 | **1.6** | 2.0 |
| Average | 18.1 | 13.8 | 7.9 | **2.7** |

### 3.5. Efficiency

Table 7 compares time efficiency. For watermark embedding, [13]'s extensive preprocessing and [14]'s text-based scaling lead to longer embedding times. In contrast, our approach embeds the watermark independently of text, resulting in shorter processing times. In terms of extraction time, our method, as well as the approaches presented in [12] and [13], exhibit similar extraction times. However, the method proposed by [14] experiences longer extraction times due to its lower resynchronization efficiency.

**Table 7**. Comparison of the time efficiency.

| Time (s) | Pramila et al. [12] | Gugelmann et al. [13] | Fang et al. [14] | Proposed |
|---|---|---|---|---|
| Embedding | 1.46 | 9.54 | 2.55 | **1.25** |
| Extraction | 6.05 | 8.33 | 26.39 | **6.04** |

## 4. CONCLUSION AND FUTURE WORK

This work introduces an enhanced screen-shooting resilient document watermark scheme that significantly enhanced imperceptibility, robustness adaptability and efficiency, encompassing these key aspects. Nevertheless, challenges remain when dealing with documents containing images and other non-textual elements, posing a potential threat to robustness. Our future work will focus on further enhancing the performance of the proposed watermarking scheme in such scenarios.

# 5. REFERENCES

[1] Haozhe Chen, Weiming Zhang, Kunlin Liu, Kejiang Chen, Han Fang, and Nenghai Yu, "Speech pattern based black-box model watermarking for automatic speech recognition," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 3059–3063.

[2] Yongbaek Cho, Changhoon Kim, Yezhou Yang, and Yi Ren, "Attributable watermarking of speech generative models," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 3069–3073.

[3] Pierre Fernandez, Alexandre Sablayrolles, Teddy Furon, Hervé Jégou, and Matthijs Douze, "Watermarking images in self-supervised latent spaces," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 3054–3058.

[4] Jack T Brassil, Steven Low, Nicholas F. Maxemchuk, and Lawrence O'Gorman, "Electronic marking and identification techniques to discourage document copying," *IEEE Journal on Selected Areas in Communications*, vol. 13, no. 8, pp. 1495–1504, 1995.

[5] Jack T Brassil, Steven Low, and Nicholas F Maxemchuk, "Copyright protection for the electronic distribution of text documents," *Proceedings of the IEEE*, vol. 87, no. 7, pp. 1181–1196, 1999.

[6] Chang Xiao, Cheng Zhang, and Changxi Zheng, "Fontcode: Embedding information in text documents using glyph perturbation," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 2, pp. 1–16, 2018.

[7] Wenfa Qi, Wei Guo, Tong Zhang, Yuxin Liu, Zongming Guo, and Xifeng Fang, "Robust authentication for paper-based text documents based on text watermarking technology," *Mathematical Biosciences and Engineering*, vol. 16, no. 4, pp. 2233–2249, 2019.

[8] Xi Yang, Weiming Zhang, Han Fang, Zehua Ma, and Nenghai Yu, "Language universal font watermarking with multi-ple cross-media robustness," *Signal Processing*, vol. 203, pp. 108791, 2023.

[9] Maikel Lázaro Pérez Gort, Martina Olliaro, Agostino Cortesi, and Claudia Feregrino Uribe, "Semantic-driven watermarking of relational textual databases," *Expert Systems with Applications*, vol. 167, pp. 114013, 2021.

[10] Jipeng Qiang, Shiyu Zhu, Yun Li, Yi Zhu, Yunhao Yuan, and Xindong Wu, "Natural language watermarking via paraphraser-based lexical substitution," *Artificial Intelligence*, vol. 317, pp. 103859, 2023.

[11] Masahiko Suzaki and Masayuki Suto, "A watermark embedding and extracting method for printed documents," *Electronics and Communications in Japan (Part III: Fundamental Electronic Science)*, vol. 88, no. 7, pp. 43–51, 2005.

[12] Anu Pramila, Anja Keskinarkaus, Valtteri Takala, and Tapio Seppänen, "Extracting watermarks from printouts captured with wide angles using computational photography," *Multimed. Tools Appl.*, vol. 76, pp. 16063–16084, 2017.

[13] David Gugelmann, David Sommer, Vincent Lenders, Markus Happe, and Laurent Vanbever, "Screen watermarking for data theft investigation and attribution," in *2018 10th International Conference on Cyber Conflict (CyCon)*. IEEE, 2018, pp. 391–408.

[14] Han Fang, Weiming Zhang, Zehua Ma, Hang Zhou, Shan Sun, Hao Cui, and Nenghai Yu, "A camera shooting resilient watermarking scheme for underpainting documents," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 11, pp. 4075–4089, 2019.

[15] Joung-Youn Kim, Lee-Sup Kim, and Seung-Ho Hwang, "An advanced contrast enhancement using partially overlapped sub-block histogram equalization," *IEEE Transactions on circuits and systems for video technology*, vol. 11, no. 4, pp. 475–484, 2001.

[16] C.X.Liu, "The three-body problem," *Chongqing Press, China*, 2008.

[17] George R. R. Martin, "A song of ice and fire," *Penguin Random House US*, 1996.