

http://www.aimspress.com/journal/MBE

Mathematical Biosciences and Engineering, 16(4): 2233–2249 DOI: 10.3934/mbe.2019110 Received: 05 January 2019 Accepted: 12 February 2019 Published: 15 March 2019

## Research article

## Robust authentication for paper-based text documents based on text watermarking technology

Wenfa Qi<sup>1</sup>, Wei Guo<sup>2</sup>, Tong Zhang<sup>1</sup>, Yuxin Liu<sup>1</sup>, Zongming Guo<sup>1,\*</sup>and Xifeng Fang<sup>3</sup>

<sup>1</sup> Institute of Computer Science and Technology, Peking University, Beijing 100871, China

<sup>2</sup> Chinese Academy of Inspection and Quarantine, Beijing, 100176, China

<sup>3</sup> Citrix System Inc, 4988 Great America Parkway Santa Clara, CA 95054, United States

\* **Correspondence:** Email: guozongming@pku.edu.cn.

**Abstract:** Aiming at the problem of easy tampering and difficult integrity authentication of paper text documents, this paper proposes a robust content authentication method for printed documents based on text watermarking scheme resisting print-and-scan attack. Firstly, an authentication watermark signal sequence related to content of text document is generated based on the Logistic chaotic map model; then, the authentication watermark signal sequence is embedded into printed paper document by using a robust text watermarking scheme; finally, the watermark information is extracted from scanned image of paper document, and compared with the authentication watermark information calculated in real time by the text document obtained by OCR technology, thereby performing content integrity authentication of paper text documents. Experimental results show that our method can achieve the robust content integrity authentication of paper text document after embedding the watermark information has a good visual effect, and the text watermarking scheme has a large information capacity.

**Keywords:** integrity authentication; chaotic mapping; image template matching; text watermarking scheme; minimum editing distance

## 1. Introduction

With the rapid development of network information technology, the communication of multimedia information has reached an unprecedented depth and breadth. The digital multimedia data provides great convenience for the wide dissemination of information. Various multimedia file data can be easily obtained through the network, and various forms of tampering can be performed on the multimedia data content. Therefore, the problem of the safe dissemination of multimedia information has become an important and urgent research topic. In recent years, how to efficiently solve the content authentication



Figure 1. The tampering identification process for paper document content.

problem of multimedia information such as images [1, 2, 3, 4, 6, 5, 7, 8, 9, 10], video [11, 12, 13, 14, 15, 16], audio [17, 18, 19, 20] and text [21, 22, 23, 24, 25] has become an increasing concern among multimedia security.

However, in daily work, many important text documents are still transmitted on paper carriers, such as government confidential documents, financial contracts and legal documents, etc. It is worth noting that paper documents can be easily falsified and forged, for example, the common methods include page-extracting, page-changing, or tampering with individual paragraphs, words, and data in the text file. Once some important documents have been tampered with, the consequences will be very serious. If some forged documents spread within a certain range, it will cause reputational damage and economic losses, while severe ones will affect production safety and even national security. Therefore, the integrity certification and authenticity identification of important paper documents is extremely critical.

In order to achieve content integrity authentication of paper documents, it is necessary to embed auxiliary authentication information in the printed paper document, digitize the paper document to obtain an image file, and extract corresponding information from the image to verify whether the document content has been changed. To this end, some researchers propose the methods to hide watermark information by printing visible shading [26, 27] or 2D barcode [28] at specific locations on each page of the paper document, which can resist print-and-scan attack. However, it may affect the normal reading of the text document with poor visual effect, which is not allowed in some cases especially for government official documents. The text watermarking technology [29, 30, 31, 32, 33, 34, 35, 36] can convert the key information of the document into invisible watermark information which can be embedded into the printed paper documents. Only after digitizing the paper documents by scanner, the hidden watermark information can be extracted through special software tools. So it can identify the authenticity of the paper documents through the extracted watermark information, thereby realizing the anti-forgery traceability, copyright protection, and content certification of the text documents and other purposes. The above text watermarking algorithms mainly focus on the robustness of resisting the print-and-copy attack, and it only needs to be able to completely extract the recovered watermark information, and does not need to ensure that each watermark information bit is correct. Therefore, this kind of method can detect the tampering behaviors of page-changing or replacement of the entire page of paper documents, but cannot detect the tampering behaviors of individual paragraphs, words and data, and accurately locate the tampering position in the whole text document.

This paper proposes a robust authentication method for the paper document content based on text

watermarking algorithm resisting print-and-scan attack. It embeds the authentication watermark signal sequence by replacing the vector fonts library to achieve the tampering identification for paper documents. The detailed content authentication process is shown in Figure 1.

#### 2. Paper document output

This section mainly discusses the output process of paper document containing the content authentication information. It generates a watermark information sequence as a paper document content authentication information according to certain rules. The information is embedded into the printed document by using text watermarking algorithm based on the replacement of vector font library. The text watermarking algorithm is capable of resisting print-and-scan attack, and it can determine whether the paper document has been tampered with and locate the corresponding tampering position based on the extracted content authentication information.

#### 2.1. Authentication Information Generation

For the sake of simplicity, security and reliability, this paper uses a Logistic chaotic map model [37] to generate watermark information bit sequence. Chaos is a deterministic and stochastic process that occurs in nonlinear dynamic systems. This process is aperiodic, non-converged but bounded, and has a sensitive dependence on the initial value [38, 39]. Using this property, the chaotic maps can provide a large number of non-correlated, stochastic, and deterministic signals that are easy to generate and regenerate.

Logistic mapping is a very simple but widely studied chaotic dynamic system, which can be described by the following nonlinear difference equation:

$$x_{n+1} = \lambda x_n (1 - x_n), \lambda \in [0, 4], x_n \in [0, 1]$$
(2.1)

The study shows that the chaotic region of Logistic chaotic map is  $\lambda \in [\lambda_{\infty}, 4]$ , where  $\lambda_{\infty} = 3.569945672\cdots$ . It has been theoretically proved that the cross-correlation of two chaotic sequences  $x_0, x_1, ..., x_n$  and  $y_0, y_1, ..., y_n$  generated by two different initial values  $x_0$  and  $y_0$  is zero, which reflects the extreme sensitivity of Logistic chaotic maps to initial values. When  $\lambda = 4$ , the probability distribution density function of the Logistic chaotic sequence is:

$$\rho(x) = \begin{cases} \frac{1}{\pi \sqrt{x(1-x)}}, & x \in (0,1) \\ 0, & others \end{cases}$$
(2.2)

the mean of the sequence is:  $\overline{x} = E\{x\} = 0.5$ . Therefore, the real-value chaotic sequence  $x_0, x_1, \dots, x_n$  can be converted into 0/1 binary sequence  $b_0, b_1, \dots, b_n$  by the threshold function  $n_0(x)$ ,

$$n_0(x) = \begin{cases} 1, & x \ge 0.5\\ 0, & x < 0.5 \end{cases}$$
(2.3)

In this paper, we perform content integrity authentication by judging whether the content of the text document output to the paper carrier has been changed or not. That is, when the characters in the text document are added, deleted or modified, it is judged that the paper document has been tampered with, which belongs to the robust certification. Since our method uses the chaotic sequence as the

authentication watermark signal, in order to fully utilize the extreme sensitivity of the chaotic system to the initial value, the construction of the authentication watermark signal sequence must be related to the original character content. If the text strings are simply grouped and the initial value calculation of the chaotic system is performed, the calculation result of the initial value will change once the character is added or deleted, and the subsequent authentication information extraction is completely wrong. Therefore, the initial value of the chaotic system must be calculated in relation to each character, and there is independence on different characters. For each character, the specific authentication watermark signal generation process is as follows:

- (1) Read the hexadecimal Unicode of character *C* occupying a total of 2 bytes;
- (2) Convert the Unicode encoding of character C into a binary bit string and divide it into four subgroups evenly, each group contains 4 bits, and four subgroups are converted into four decimal numbers  $\sigma_1$ ,  $\sigma_2$ ,  $\sigma_3$  and  $\sigma_4$ , respectively;
- (3) For the character *C*, calculate the initial value  $x_0$  of the Logistic chaotic map as follows:

$$x_0 = \frac{1}{2} \left( \frac{1}{4} \sum_{i=1}^4 \frac{\sigma_i}{16} + k \right)$$
(2.4)

where k is the watermark key,  $k \in (0, 1)$ .

- (4) Take  $x_0$  and  $\lambda = 4$  as the initial value and the parameter of the Logistic chaotic map, and substitute the Logistic map iteration equation (2.1), and generate a real-value chaotic sequence  $x_0, x_1, \dots, x_{15}$  of length 16 which is converted into a binary sequence using the threshold function  $n_0(x)$ , and take the 4th bit of the binary sequence as the watermark bit corresponding to the character *C*;
- (5) Each character is calculated according to the above process to obtain a corresponding watermark information bit string sequence  $\varphi_1 = \{b_0, b_1, ..., b_m\}$ , where *m* is the number of all characters;
- (6) Calculate the CRC16 (Cyclic Redundancy Check) bit string of the watermark information bit string sequence  $\varphi_1$ ;
- (7) Finally, the CRC16 check code is appended to the watermark information bit sequence  $\varphi_1$  to generate a final authentication watermark signal sequence  $\varphi_2$ .

## 2.2. Certification Information Embedding

We use the text watermarking algorithm based on vector font replacement to embed the authentication watermark signal sequence  $\varphi_2$  calculated above into the printed paper document. The specific process is as follows: Firstly, a special watermark vector font library is designed, that is, by appropriately changing the topological structure of characters or strings, multiple glyphs of semantically identical characters (strings) are designed, and each glyph is given a different character code representing different watermark information bits, and save all character codes to a new font library. For example, as shown in Figure 2, the position of character stroke marked with gray shifts horizontally; then, when the text document is printed out, the text document content is intercepted by the terminal monitoring service program in the computer; finally, according to the watermark information bit string to be embedded, the watermark information is embedded by dynamically replacing the different glyph



**Figure 2.** Schematic diagram of watermark character structure design. (a) The original FangSong font Chinese character "在". (b) The character structure is obtained by shifting the gray stroke to right direction and represents the watermark information bit 0. (c) The character structure is obtained by shifting the gray stroke to left direction and represents the watermark information bit 1.



**Figure 3.** Glyph structures morphing of Arabic numeral 5: the left one is a standard font, the middle one represents the watermark information bit 0 with the gray stroke shifting downward, and the right one represents the watermark information bit 1 with the gray stroke shifting upward.

structures of the characters, and the modified text document content data is sent to the physical printer to be output in the end.

It should be noted that this method can be applied not only to complex Chinese characters, but also to English letters and Arabic numerals, just as shown in Figure 3 and Figure 4, the position of character strokes marked with gray shifts vertically. Considering the information capacity, each character can map multiple different glyph structures to represent more watermark information. For example, when mapping to 4 different glyphs, each character can represent a 2-bit watermark information. Therefore, the capacity of the authentication watermark signal sequence  $\varphi_2$  calculated by all the characters is only close to half of the watermark capacity. Other extended parameters information including the watermark key k, may also be encoded and embedded into the printed paper document. If the number of characters in the text document is large, the authentication watermark signal sequence  $\varphi_2$  will be repeatedly embedded to resist the watermark information missing attack caused by the multi-character deletion. In addition, all the characters in the printed document after embedding the watermark information have been replaced with the corresponding characters in the watermark font library.

# FFF

**Figure 4.** Glyph structures morphing of English character F: the left one is a standard font, the middle one represents the watermark information bit 0 with the gray stroke shifting upward, and the right one represents the watermark information bit 1 with the gray stroke shifting downward.

#### 3. Paper document content certification

The digital document image *I* obtained by scanning the paper document with the watermark information embedded is processed, and the difference between the watermark character glyph structures is detected by the normalized-character-image-template-match method, thereby extracting the previously embedded authentication watermark signal sequence  $\varphi'_2$  and recovering the original watermark information bit sequence  $\varphi'_1$ ; and it recalculates the watermark information bit sequence  $\varphi''_1$  according to the characters of the text document image *I*; finally by comparing the difference between the signal sequences  $\varphi'_1$  and  $\varphi''_1$ , it determines whether the content of the paper document has been tampered with.

#### 3.1. Normalized Character Image Template Matching

Let the template character image be A of size  $M_A * M_B$ , the character image to be matched is B, and the size is  $N_A * N_B$ , where  $M_A \le N_A$ ,  $M_B \le N_B$ . The matching process of the character image is to calculate the similarity of the two images, that is, calculate the autocorrelation coefficients of the character image A and the character image B at different positions, and the maximum of all the coefficients is the similarity of the two images.

As shown in Figure 5, when the template image A slides to (u, v) on the image B to be matched, the normalized correlation coefficient r(u, v) between the two images is calculated as follows [40]:

$$r(u,v) = \frac{\sum_{i,j} [f(u+i,v+j) - \overline{f}_{u,v}][t(i,j) - \overline{t}]]}{\sqrt{(\sum_{i,j} [f(u+i,v+j) - \overline{f}_{u,v}]^2)(\sum_{i,j} [t(i,j) - \overline{t}]^2)}}$$
(3.1)

Where f(i, j) is mean value of the character image *B* at the (i, j) position, and t(i, j) is the mean value of the template character image *A* at the (i, j) position,  $\bar{t}$  is pixel mean value of the template image *A*,  $\bar{f}$  is the pixel mean value of the corresponding image block to be matched under the position where the current template image *A* is located.

Due to high computational complexity of above formula (3.1), the equivalent calculation method is as follows :

Firstly, let  $t'(i, j) = t(i, j) - \overline{t}$ , then the numerator of equation (3.1) can be expressed as follows:

$$\sum_{i,j} f(u+i,v+j)t'(i,j) - \bar{f} \sum_{i,j} t'(i,j)$$
(3.2)

Mathematical Biosciences and Engineering

Volume 16, Issue 4, 2233–2249



Figure 5. Character image matching diagram.

Obviously, the latter term of equation (3.2) is always equal to zero, and the previous term of equation (3.2) can be obtained by:

$$\sum_{i,j} f(u+i,v+j)t(i,j) - \bar{t} \sum_{i,j} f(u+i,v+j)$$
(3.3)

The first term of the equation (3.3) can be regarded as a convolution operation of two signals in the spatial domain which is equivalent to the frequency multiplication operation based on the frequency domain. Therefore, the first term of equation (3.3) is equivalent to the following form:

$$F^{-1}[F(f)F^*(t)]$$
(3.4)

where, F is the Fourier transform of the original signal,  $F^*$  is the conjugate complex operation of the transformed result, and  $F^{-1}$  is the inverse Fourier transform for the frequency domain signal.

Next, for the denominator portion of the equation (3.1), the latter term is the pixel value variance of the template character image A. The first item is simplified to get the following results:

$$\sum_{i,j} \left[ f(u+i,v+j) \right]^2 - \frac{\left[ \sum_{i,j} f(u+i,v+j) \right]^2}{M_A * N_A}$$
(3.5)

In order to calculate the equation (3.5), the sum and square sum of all pixel values in the sliding window of the character image B at (u, v) should be obtained firstly. To do this, the following functions are defined:

$$s(u, v) = f(u, v) + s(u - 1, v) + s(u, v - 1) - s(u - 1, v - 1)$$
(3.6)

$$s_2(u,v) = f^2(u,v) + s_2(u-1,v) + s_2(u,v-1) - s_2(u-1,v-1)$$
(3.7)

Then, through the dynamic programming, the accumulation of all elements of the sub-matrix composed of the point (0, 0) to the point (u, v) as the diagonal element in the character image *B* is calculated by looking up the table, and the squared sum is obtained by summing the squares of all the elements in the submatrix. After that, the cumulative sum and cumulative square sum of all pixels in the windows which locates at any position (u, v) of the size  $(M_A, N_A)$  can be calculated as follows:

$$\sum_{i,j} f(u+i,v+j) = s(u+M_A-1,v+N_A-1) - s(u-1,v+N_A-1) - s(u+M_A-1,v-1) + s(u-1,v-1)$$
(3.8)

$$\sum_{i,j} f^2(u+i,v+j) = s_2(u+M_A-1,v+N_A-1) - s_2(u-1,v+N_A-1) - s_2(u+M_A-1,v-1) + s_2(u-1,v-1)$$
(3.9)

As shown in Figure 5, when the template character image A slides to the point (u, v), the cumulative sum and cumulative square sum of pixels in the sliding window, is equal to the corresponding values in the sub-matrix portion marked by yellow in the character image B. By substituting equations (3.8) and (3.6) into equations (3.3) and (3.5), it is possible to find the normalized correlation coefficient when the template character image A slides along the character image B to the point (u, v). The matching result of image A and B is the maximum value among all the correlation coefficients of the two graphs. The position coordinate  $(u_m, v_m)$  at this time is the offset point coordinate of the image A with respect to the character image B, that is, when the image A moves to  $(u_m, v_m)$  and the matching degree of the image B is the highest.

#### 3.2. Watermark information extraction

For each character in the text document image *I*, it matches the scanned character image with the standard glyph image of the character and each variant glyph image one by one by using the normalized image template matching method described in Section 3.1, and makes the following judgments:

- a) If the matching degree of scanned character image and standard glyph image is the highest, it indicates that the character is not the character after embedding the watermark in the printout process, and does not represent any watermark information. Therefore, the character can be regarded as a new character that is falsified after being added and inserted into the contents of the paper document, and the position of the falsified character is recorded;
- b) Otherwise, according to the matching degree of scanned character image and the variant glyph image, the binary watermark information string represented by the character is calculated.

After all the characters are subjected to the image matching operation, all the watermark information bit strings are recovered and the final authenticated watermark signal sequence  $\varphi'_2$  is obtained.

#### 3.3. Document Content Robust Authentication

The process of the robust authentication of the paper document content is described as follows:

- a) Decompose the obtained authentication watermark signal sequence  $\varphi'_2$  into two parts which are the embedded watermark information bit string sequence  $\varphi'_1$  and CRC16 check code bit string information respectively;
- b) For each character in the text document image *I*, recalculate the watermark information bit string sequence  $\varphi_1''$  according to the method described in Section 2.1;



Figure 6. The sub-image block of Chinese character "在" after printing and scanning.

- c) Perform the same CRC check operation for the extracted watermark information bit sequence  $\varphi'_1$ , if the check passes, the watermark information bit sequence  $\varphi'_1$  is correct, and the watermark signal sequences  $\varphi'_1$  and  $\varphi''_1$  are compared one bit by one bit: If the two information bit strings are the same, the content of the paper document does not change, otherwise it is falsified, and the accurate tampering position is given according to the difference of the bit string;
- d) If the CRC check fails, the similarity comparison between the signal sequence  $\varphi'_1$  and  $\varphi''_1$  is executed, and the longest common substring of the two signal sequences is calculated according to the string minimum edit distance algorithm [41]. The positions of the rest different string bits are regarded as the suspicious tampering positions.

### 4. Experimental results and discussion

At first, we verify the robustness of the text watermarking algorithm based on the vector font library. After modifying the glyph topology of the original Chinese character " $\pm$ " as shown in Figure 2(a), it is mapped into two different glyph structures, as shown in Figure 2(b) and Figure 2(c), respectively, and they represent the watermark information bit 0 and 1, respectively. When the watermark information is embedded, according to the watermark bit string information, Figure 2(a) is replaced with the two glyph structures of Figure 2(b) or Figure 2(c) correspondingly. When the watermark information is extracted, the paper document embedded with the watermark information is scanned into a digitized text image, and the sub-image block of the Chinese character " $\pm$ " is segmented as shown in Figure 6.

The normalized-character-image-template-match (NCITM) method described in Section 3.1 is used to match the glyph structure images of Figure 6 with the two images Figure 2(b) and Figure 2(c), respectively. It should be noted that before image matching, the size of the scanned character image needs to be normalized. Firstly, the paper document is scanned with the accuracy of more than 300 PDI to ensure that the scanned character image is clear; then, the size of each character is scaled to the size of the standard character image read from the font library; finally, the NCITM process is carried out. The matching results are shown in Figure 7, in which Figure 7(a) and Figure 7(b) are the schematic diagrams showing the optimal matching effect of Figure 6 and the two structures shown in Figure 2(b) and Figure 2(c), respectively. The correlation coefficients of the glyph structure calculated according



**Figure 7.** The matching result of Figure 6 with the standard watermark glyph structure. (a) The optimal matching effect of Figure 6 and Figure 2(b). (b) The optimal matching effect of Figure 6 and Figure 2(c).

to formula (3.1) are 0.94249 and 0.79489, respectively. It can be decided that the matching degree of Figure 6 and Figure 2(b) is higher, thereby deriving the watermark information bit string represented by the character as 0. In order to improve the information capacity, each character can be mapped to more deformation effects, for example, when one character corresponds to 4 different deformations, and one character can represent 2-bit watermark information. However, when the number of corresponding deformed fonts of one character increases, the workload of designing the watermark font library is relatively larger.

In this experiment, we designed the robust performance test of the text watermarking algorithm in three scenarios: printing and scanning, capturing screen and shooting screen. Figure 8 is a schematic diagram of the effect of the original document without embedding watermark information with the Chinese FangSong font type of size 13 pounds. Figure 9 shows the effect of the document after embedding the watermark information, and each character has different font deformation effects. For example, the Chinese character "不" has 4 different deformations and represents 2 bits watermark information, while the Chinese character "人" has two different deformations and represents one bit information. In order to facilitate the description of the font deformation effect, the Chinese characters "不" and "人" are enlarged and bolded intentionally. Figure 9 contains 168 Chinese characters in which 326 bits watermark information is embedded. Figure 10 is the image data obtained by scanning the paper document shown in Figure 9, and this test is marked as PS. Figure 11 is the document image effect after the screen shooting, while the electronic document format with the watermark information embedded is displaying on the computer screen, which is marked as CC. Figure 12 is a schematic diagram of photographing computer screen of Figure 10 outside the screen through the mobile phone, and this test is marked as PC. It can be seen from the Figure 12 that there is serious moiré interference phenomena in the photograph. In addition, for the test of shooting screen of the electronic document just like Figure 11, we use MS Office Word software to open the doc/docx format file embedded with watermark information under the different page displaying ratios of 100%, 110%, 120%, 130%, 140% and 150% respectively, then a screen capture is taken and the screen capture image is obtained from which the watermark information is extracted. Under the attack operations shown in Figure 10-12 as mentioned

缓歌慢舞凝丝竹,尽日君王看不足。渔阳鼙鼓动地来,惊破霓裳羽衣曲。 九重城阙烟尘生,千乘万骑西南行。翠华摇摇行复止,西出都门百余里。 六军不发无奈何,宛转蛾眉马前死。花钿委地无人收,翠翘金雀玉搔头。 君王掩面救不得,回看血泪相和流。黄埃散漫风萧索,云栈萦纡登剑阁。 峨嵋山下少人行,旌旗无光日色薄。蜀江水碧蜀山青,圣主朝朝暮暮情。 行宫见月伤心色,夜雨闻铃肠断声。天旋地转回龙驭,到此踌躇不能去。 Figure 8. The original document image.

缓歌慢舞凝丝竹,尽日君王看不足。渔阳鼙鼓动地来,惊破霓裳羽衣曲。 九重城阙烟尘生,千乘万骑西南行。翠华摇摇行复止,西出都门百余里。 六军不发无奈何,宛转蛾眉马前死。花钿委地无人收,翠翘金雀玉搔头。 君王掩面救不得,回看血泪相和流。黄埃散漫风萧索,云栈萦纡登剑阁。 峨嵋山下少人行,旌旗无光日色薄。蜀江水碧蜀山青,圣主朝朝暮暮情。 行宫见月伤心色,夜雨闻铃肠断声。天旋地转回龙驭,到此踌躇不能去。 Figure 9. The effect of the original watermarked document.

above, the success rate of watermark extraction of [42], [43] and our proposed method are shown in Table 1. It can be seen that our method achieves very good extraction success rate except for some special cases where the display ratios are 100% and 110%, because the deformations of the character strokes in the screenshot image are severe. When the display ratio of the electronic document page is larger, the clearer the stroke structure of the character is displayed and the higher the accuracy of the character image is, and correspondingly the higher the watermark extraction rate is. So we can conclude that the text watermarking algorithm has strong robust performance for attack operations such as print-and scan, capturing screen and shooting screen, etc. By contrast, our approach has outperformed [42] and [43] in all cases.

As mentioned above, in the text-based watermarking algorithm based on vector fonts, we change the glyph structure of the watermark characters, which will remain unchanged before and after the

	PS	<i>CC</i> -100%	<i>CC</i> -110%	<i>CC</i> -120%	<i>CC</i> -130%	<i>CC</i> -140%	<i>CC</i> -150%	PC
Method[42]	100%	51%	79%	86%	88%	93%	95%	94%
Method[43]	100%	75%	85%	88%	92%	96%	98%	100%
Our method	100%	89%	97%	100%	100%	100%	100%	100%

 Table 1. The watermark correct extraction rate.

Mathematical Biosciences and Engineering

缓歌慢舞凝丝竹,尽日君王看不足。渔阳鼙鼓动地来,惊破霓裳羽衣曲。 九重城阙烟尘生,千乘万骑西南行。翠华摇摇行复止,西出都门百余里。 六军不发无奈何,宛转蛾眉马前死。花钿委地无人收,翠翘金雀玉搔头。 君王掩面救不得,回看血泪相和流。黄埃散漫风萧索,云栈萦纡登剑阁。 峨嵋山下少人行,旌旗无光日色薄。蜀江水碧蜀山青,圣主朝朝暮暮情。 行宫见月伤心色,夜雨闻铃肠断声。天旋地转回龙驭,到此踌躇不能去。

Figure 10. The effect of the printed and scanned watermarked paper document.

缓歌慢舞凝丝竹,尽日君王看不足。渔阳鼙鼓动地来,惊破霓裳羽衣曲。 九重城阙烟尘生,千乘万骑西南行。翠华摇摇行复止,西出都门百余里。 六军不发无奈何,宛转蛾眉马前死。花钿委地无人收,翠翘金雀玉搔头。 君王掩面救不得,回看血泪相和流。黄埃散漫风萧索,云栈萦纡登剑阁。 峨嵋山下少人行,旌旗无光日色薄。蜀江水碧蜀山青,圣主朝朝暮暮情。 行官见月伤心色,夜雨闻铃肠断声。天旋地转回龙驭,到此踌躇不能去。

Figure 11. The effect of the screen capture image while the electronic watermarked document is displaying.

缓歌慢舞凝丝竹,尽日君王看不足。渔阳鼙鼓动地来,惊破霓裳羽衣曲。 九重城阙烟尘生,千乘万骑西南行。翠华摇摇行复止,西出都门百余里。 六军不发无奈何,宛转蛾眉马前死。花钿委地无人收,翠翘金雀玉搔头。 君王掩面救不得,回看血泪相和流。黄埃散漫风萧索,云栈萦纡登剑阁。 峨嵋山下少人行,旌旗无光日色薄。蜀江水碧蜀山青,圣主朝朝暮暮情。 行宫见月伤心色,夜雨闻铃肠断声。天旋地转回龙驭,到此踌躇不能去。

**Figure 12.** The effect of photographing computer screen through the mobile phone while the electronic watermarked document is displaying.

print-and-scan. Therefore, the algorithm can effectively resist the print-and-scan attack, the watermark information is propagated along with the paper document, and if necessary, the paper document can be digitized into an image file, and the content integrity authentication of the paper document can be performed by extracting the hidden watermark information.

In this experiment, the authentication watermark signal sequence is calculated according to the document content as shown in Figure 13 using the method described in Section 2.2, and the watermark information is embedded in the process of document printing and output. Figure 14 is a schematic diagram showing the effect of manually tampering with the contents of a paper document. The watermark information extraction is performed according to the NCITM method described in Section 3.1. The characters "红", "女", "0", and "监" identified in the rectangle are new characters added by deleting the original ones in the same positions, which does not disturb the character order in the original sequence and is seen as the Character-Replacement mode. These characters have the highest matching degrees with the characters in the standard font library and do not contain any watermark information, and they do not belong to the watermark fonts. Since each character in Figure 13 has been modified by font substitution, the above four characters are identified as content tampering and the corresponding tampering positions are recorded. Similarly, the Arabic numerals 1 marked with hexagons is also confirmed as being tampered with by the Character-Adding mode, and the corresponding tampering position is recorded. Finally, using the document content robustness authentication process described in Section 3.3, the similarity between the embedded watermark information bit sequence  $\varphi'_1$  and the recalculated watermark information bit sequence  $\varphi_1''$  is compared. The part between the Chinese character "[" and "]" will be marked as content tampering by Character-Deleting mode, and the corresponding tampering position will be recorded.

#### 5. Conclusion

This paper proposes a robust content authentication method for paper text documents, which solves the problem of paper document content integrity verification, tamper identification and tamper position locating, etc. The main contributions include: 1) A Logistic chaotic map model is used to generate a watermark information bit sequence, which is related to the content of the text document. Based on the sensitivity of the Logistic chaotic map to the initial value, once the corresponding character changes, the watermark information bit string will be wrong, thereby performing the content tampering identification and locating the specific tampering position. 2) The generated authentication watermark signal sequence is embedded into the printed paper document by using the text watermarking algorithm based on the vector font library. The watermark information does not depend on the additional carrier, and the text document embedded with the watermark has a good visual effect. The watermark information is embedded during the process of printout, which avoids the risk of malicious tampering of the watermark information, so this method has high security performance. The watermarking algorithm is robust against print-and-scan attack and has a large information capacity. 3) The minimum string edit distance algorithm is used to compare the difference between authentication watermark signal sequence extracted from the scanned image and the one calculated in real-time, through which it discriminates whether the paper text document content changes and accurately locates the tampering position. In addition, there are still some issues to be improved for our method, such as how to establish the association between the text document content and the logistic chaotic map based on the logistic chaotic

# 在职证明(Sample)

兹证明<u>黄飞鸿</u>同志,汉语拼音<u>Huang Feihong</u>,性 别(Gender)<u>男</u>,民族(Nationality)<u>汉</u>,身份证号码(ID Number)<u>123456789012345678</u>,自<u>2013年</u>10月20日至 <u>2018年3月2</u>日在我单位工作,担任本公司<u>技术总经理</u> 一职,月收入约为<u>5000</u>元。

公司盖章: Love future technology Limited 公司负责人签名: Tommy

日期: 2018.06.10

Figure 13. The original printed paper document.

在职证明(Sample)							
兹证明 <u>黄飞红</u> 同志,汉语拼音 <u>Huang Feihong</u> ,性							
别(Gender) <u>安</u> ,民族(Nationality) <u>汉</u> ,身份证号码(ID							
Number) <u>123456789012345678</u> ,自 <u>2010</u> 年 <u>10</u> 月 <u>20</u> 日至							
2018 年 3 月 2 日在我单位工作,担任本公司 技术总监							
一职,月收入约为 <u>05000</u> 元。							
公司盖章: Love future technology Limited							
公司负责人签名: Tommy							
日期: 2018.06.10							

Figure 14. The manually falsified paper documen of Figure 10.

2246

map model more reasonably, and to locate the tampering position more accurately.

## Acknowledgments

This work was supported by National R&D project of China under contract No. 2018YFB0803402 and National Natural Science Foundation of China (Nos.U1636206,61572052).

## **Conflict of interest**

The authors declare no competing interests.

## References

- 1. J. Lee and C. S. Won, A watermarking sequence using parities of error control coding for image authentication and correction, *IEEE Trans. Consum. Electron.*, **46** (2000), 313–317.
- 2. C. Li, Y. Wang and B. Ma, et al., Multi-block dependency based fragile watermarking scheme for fingerprint images protection, *Multimed. Tools Appl.*, **64** (2013), 757–776.
- 3. X. Lv and Z. J. Wang, Perceptual image hashing based on shape contexts and local feature points, *IEEE Trans. Inf. Forensics Security*, **7** (2012), 1081–1093.
- 4. K. Maeno, Q. Sun and S. F. Chang, et al., New semi-fragile image authentication watermarking techniques using random bias and nonuniform quantization, *IEEE Trans. Multimedia*, **8** (2006), 32–45.
- 5. C. Qin, P. Ji and C. Chang, et al., Non-uniform Watermark Sharing Based on Optimal Iterative BTC for Image Tampering Recovery, *IEEE Multimedia*, **25** (2018), 36–48.
- 6. Z. Tang, X. Zhang and L. Huang, et al., Robust image hashing using ring-based entropies, *Signal Process.*, **93** (2013), 2061–2069.
- 7. R. Vartak and S. Deshmukh, Survey of digital image authentication techniques, *International Journal of Research in Advent Technology*, **2** (2014), 176–179.
- 8. Y. Zhao, S. Wang and X. Zhang, et al., Robust hashing for image authentication using zernike moments and local features, *IEEE Trans. Inf. Forensics Security*, **8** (2013), 55–63.
- 9. Z. Yang, Y. Huang and X. Li, et al., Efficient secure data provenance scheme in multimedia outsourcing and sharing, *CMC-Comput. Mat. Contin.*, **56** (2018), 1–17.
- 10. C. Qin, X. Chen and X. Luo, et al., Perceptual Image Hashing via Dual-cross Pattern Encoding and Salient Structure Detection, *Inf. Sci.*, **423** (2018), 284–302.
- 11. J. Dittmann, A. Steinmetz and R. Steinmetz, Content-based digital signature for motion pictures authentication and content-fragile watermarking, in *Proceedings IEEE International Conference* on Multimedia Computing and Systems, IEEE, (1999), 209–213.
- 12. Z. Hou and X. Tang, Integrity authentication scheme of color video based on the fragile watermarking, in 2011 International Conference on Electronics, Communications and Control (ICECC), IEEE, (2011), 4354–4358.

2247

- 13. T. S. Rao and R. R. Kurra, A smart intelligent way of video authentication using classification and decomposition of watermarking methods, preprint, arXiv:1404.7237.
- 14. D. Xu, R. Wang and J. Wang, A novel watermarking scheme for H.264/AVC video authentication, *Sig. Proc. Image Comm.*, **26** (2011), 267–279.
- 15. W. Zhang, *Research on algorithm of digital video watermarking based on H.264 for copyright protection and content*, Ph.D thesis, Beijing university of Post and telecommunications, 2013.
- 16. H. Zhu, *The research of watermarking technology based on H.264/AVC*, Ph.D thesis, Ningbo University, 2011.
- 17. J. Li, Digital watermarking technology for protecting the copyright of audio aggregation and authenticating the integrity of audio aggregation, Ph.D thesis, Ningbo University, 2012.
- 18. R. Wang and Y. Xiong, A novel watermarking algorithm for protecting audio aggregation based on ICA, in *6th International Conference on Digital Content, Multimedia Technology and its Applications*, IEEE, (2010), 302–308.
- 19. X. Wang, P. Niu and M. Lu, A robust digital audio watermarking scheme using wavelet moment invariance, *J. Syst. Softw.*, **84** (2011), 1408–1421.
- 20. I. K. Yeo and H. J. Kim, Modified patchwork algorithm: a novel audio watermarking scheme, *IEEE Trans. Speech and Audio Processing*, **11** (2003), 381–386.
- Z. Jalil, A. M. Mirza and H. Jabeen, Word length based zero-watermarking algorithm for tamper detection in text documents, in 2010 2nd International Conference on Computer Engineering and Technology, IEEE, (2010), 376–378.
- 22. Q. C. Li and Z. H. Dong, Novel text watermarking algorithm based on chinese characters structure, in 2008 International Symposium on Computer Science and Computational Technology, IEEE, (2008), 348–351.
- 23. J. Liu, L. He and D. Fang, et al., A text digital watermarking for chinese word document, in 2008 *International Symposium on Computer Science and Computational Technology*, IEEE, (2008), 217–220.
- 24. L. Xiang, Y. Li and W. Hao, et al., Reversible natural language watermarking using synonym substitution and arithmetic coding, *CMC-Comput. Mat. Contin.*, **55** (2018), 541–559.
- 25. L. Xiang, W. Wu and X. Li, et al., A linguistic steganography based on word indexing compression and candidate selection, *Multimed. Tools Appl.*, **77** (2018), 28969–28989.
- 26. W. Guo, Y. Liu and B. Yang, et al., Research on information hiding technology in paper-based document, *China Print. Pack. Study*, **5** (2013), 30–34.
- 27. B. Yang, W. Shi and W. Qi, et al., Methods and apparatus for embedding and detecting digital watermarks in a text document, 2012, US Patent 8,107,129.
- 28. C. Ji, Y. Song and Z. Pang, A tamper proof method for the paper documents, 2012, China Patent 102,722,737.
- 29. T. Amano and D. Misaki, A feature calibration method for watermarking of document images, in *Proceedings of the Fifth International Conference on Document Analysis and Recognition. IC-DAR'99 (Cat. No. PR00318)*, IEEE, (1999), 91–94.

- 30. R. Chen, *Research of Digital Watermarking technology to Resist Printing and Scanning*, Ph.D thesis, Xidian University, 2014.
- 31. M. Lei, Y. Yang and R. Hu, A print resilient watermark scheme based on character complexity, *J. B.J. Univer. Posts Telecom.*, **38** (2015), 58–62.
- 32. S. Li, A Study of Digital Watermarking Algorithm On Binary Text Image to Resist Hard-copy Attack, Ph.D thesis, Xidian University, 2012.
- 33. W. Qi, X. Li and B. Yang, et al., Document watermarking scheme for information tracking, *J. Communicat.*, **29** (2008), 183–190.
- 34. M. Wu and B. Liu, Data hiding in binary image for authentication and annotation, *IEEE Trans. Multimedia*, **6** (2004), 528–538.
- 35. H. Yang and A. C. Kot, Pattern-based data hiding for binary image authentication by connectivitypreserving, *IEEE Trans. Multimedia*, **9** (2007), 475–486.
- 36. L. Xiang, J. Yu and C. Yang, et al., A word-embedding-based steganalysis method for linguistic steganography via synonym-substitution, *IEEE Access*, **6** (2018), 64131–64141.
- 37. S. Chen and H. Leung, Ergodic chaotic parameter modulation with application to digital image watermarking, *IEEE Trans. Image Process.*, **14** (2005), 1590–1602.
- 38. Z. Li, *Image authentication based on Digital Watermarking and Its Key Issues*, Ph.D thesis, Beijing Jiaotong University, 2008.
- 39. G. Shao and M. Zhang, Novel fragile authentication watermark based on chaotic system, in 2004 *IEEE International Symposium on Industrial Electronics*, IEEE, (2004), 1491–1494.
- 40. J. C. Yoo and T. H. Han, Fast normalized cross-correlation, *Circuits Syst. Signal Process.*, **28** (2009), 819–843.
- 41. Y. Niu and C. Zhang, Comparison of string similarity algorithm, *Comput. Digit. Eng.*, **40** (2012), 14–17.
- 42. D. Liu, *Research on novel text digital watermarking technologies and their typical applications*, Ph.D thesis, University of Electronic Science and Technology, 2007.
- 43. Y. Liu, W. Guo and W. Qi, Researches on Text Image Watermarking Scheme Based on the Structure of Character Glyph, *App. Mechan. Material.*, **731** (2015), 1163–1168.



 $\bigcirc$  2019 the Author(s), licensee AIMS Press. This is an open access article distributed under the terms of the Creative Commons Attribution License (http://creativecommons.org/licenses/by/4.0)