

Forgery Detection and Localization in Scanned Documents

Husen Issa¹, Saleh Ibrahim², Soubhi Hadri³ and Yusra Abdulrahman⁴¹ Faculty of Mechanical and Electrical Engineering, Damascus University, Syria² Faculty of Engineering, University of Debrecen, Hungary³ Software Engineer, Microsoft Corporation, USA⁴ Faculty of Informatics Engineering, University of Aleppo, Syria

IJASR 2020

VOLUME 3

ISSUE 5 SEPTEMBER – OCTOBER

ISSN: 2581-7876

Abstract – Our reliance on the scanned documents is no longer a deniable fact. However, with the spread of forgery cases, their credibility started to be questionable. Recently, an extensive research has been carried out in the field of digital images forgery detection, yet only a small portion specialized in the forgery detection in scanned documents images. The distinct properties of documents (e.g. the similar-looking characters and white identical background) impose new level of challenges for forgery detection. In this paper, we adapt the identical copy-move detection, error level analysis (ELA), and metadata analysis methods for the detection and localization of specific forgery cases in scanned documents images. We use the identical copy-move approach and error level analysis to localize duplicated blocks and inserted texts within documents respectively, while we exploit the image's metadata to track the manipulation traces. Eventually, various experiments were conducted and multiple metrics were used for the assessment of our contribution. Our results show promising potentials for these algorithms to be deployed in the field of scanned documents forgery detection.

Keywords: scanned documents, forgery detection, copy-move detection, error level analysis, metadata analysis.

1. Introduction

There is no doubt that we are living in the visual information era, where images occupy the lion's share of this transformation with their persuasive and aesthetic power [1]. With the quantum leap of digitization, digital images take the lead in terms of prevalence. It only takes a glance to realize the reliance on digital images in court of laws, news, financial documents [2], and recently social media. The increasing discovered forgery cases affect negatively the authenticity and credibility of digital images, and this only becomes more challenging with the availability of powerful image editing programs that facilitate images forgery [3]. All of this evoked an extensive research in the field of digital images forgery detection, figure (1), [4], yet and despite its prominence, only a small portion handled the forgery detection of scanned text documents images. The distinct properties of scanned text documents such as the white identical background and similar-looking characters can lead to poor performance for a plain forgery detection algorithm, like most of copy-move forgery detection algorithms [5].

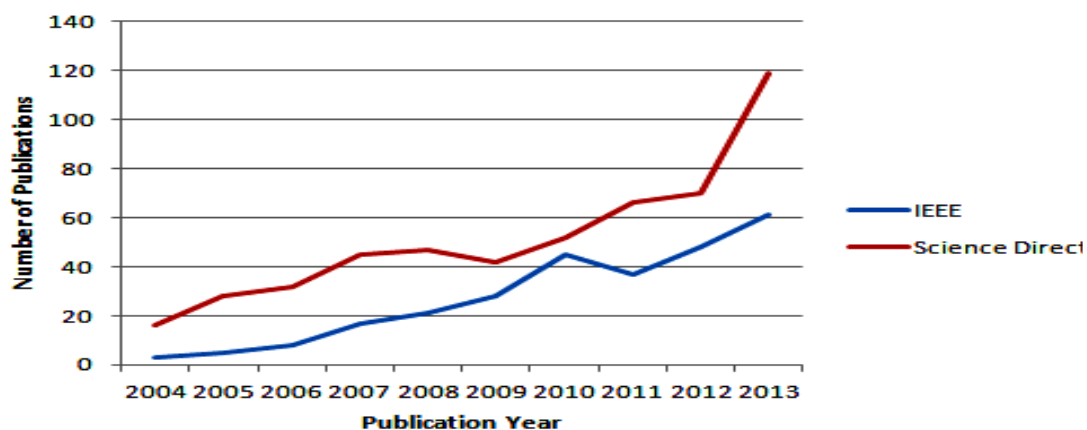


Figure (1): the number of publications in digital image forgery detection in IEEE and Science Direct.

In the domain of forgery detection in digital images, metadata information was used by many researchers as a primary tool. [6], for example, analysed the decoding properties of EXIF Metadata to detect the forged images. [7] used EXIF headers and JPEG quantization tables to extract the camera signature and compare it afterwards with a database of known authentic camera signatures. The camera signature consists of seven components. Three components are extracted from the image (image dimensions, quantization table, and Huffman code). The other three are extracted from the thumbnail (thumbnail dimensions, quantization values, and Huffman codes). While the last component comes from the EXIF metadata (entry counts from the standard IFDs, additional IFDs, entries in these additional IFDs, parser errors). Error Level Analysis(ELA) is another method that had been used in this application. [8] combined ELA with wavelet soft thresholding to clearly show the ELA results by removing the noisy components. First, the Daubechies wavelet transform is applied on each of the image channels, and after then, the noisy components are omitted using a calculated threshold for each of these channels. On another side, some researchers deployed deep learning methods in digital images forgery detection. [9] used ELA with deep learning to build a robust system which can detect forged images, where the output of ELA is an input for the convolutional neural network. [10] built a multi-stream version of Faster R-CNN model to detect image tampering, where they suggested improving accuracy of their method by adding a second stream to the model whose input is the ELA output.

Speaking of the research into forgery detection in the images of scanned documents, the relatively limited number of publications is well noticed when comparing it with the one of forgery detection in general images.[11] used block-based copy-move detection algorithm using Zernike and Hu moments. First off, an Optical Character Recognition (OCR) system is used to remove the blank spaces and false black areas resulted from the scanning. Then Zernike and Hu moments are extracted from blocks and used as features to detect forgery cases. Eventually, the results are evaluated for each moment separately. [12] Depended on that, in most forgery cases, the manipulated area in the document presents texture discontinuities with respect to other areas and more specifically in the background area around the forged text. Their method consists of dividing the image into patches, calculating the Local Binary Patterns descriptor for each patch, and finally classifying each patch in comparison to its neighbours as forged or non-forged patch using a SVM classifier. [5] Also used a copy-move detection algorithm to detect the duplicated regions in the scanned documents. Numerous experiments were carried out in their research to evaluate the performance of general copy-move forgery detection algorithms when handling scanned documents. OCR system was used in their study to specify the block size, while some image pre-processing operations were performed.

2. Proposed System

In this paper weproposean end-to-end system, figure (2) that adapt and utilize three known methods to detect and localize common forgery cases in scanned documents images. We use block-based copy-move detection and error level analysis methods to detect and localize the identical duplications and inserted black coloured texts respectively within documents images. Additionally, we extract useful information form the image metadata that might help in detecting the manipulated documents.

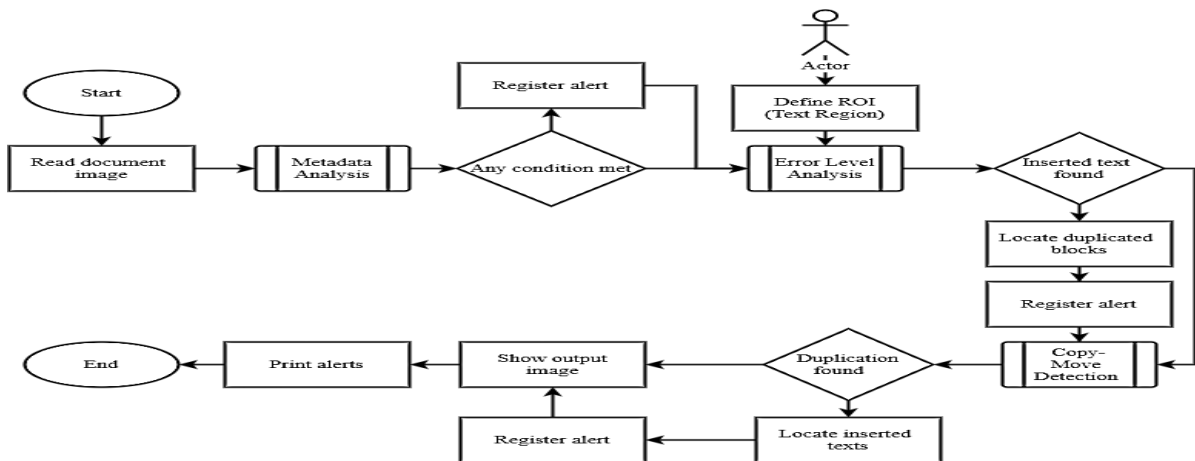


Figure (2): design of the proposed system.

2.1 Metadata Analysis

In the search for a general method that can be used in digital forensics to confirm the image manipulation, the image metadata might be the answer where other methods fail [6]. The image metadata provides detailed information on how the image was created and processed, which can be used to track the history of the image. Unfortunately, even this information can be edited or deleted [13] with specific tools. For that reason, this method is based on the assumption of non-edited metadata by the manipulator. Four types of metadata were used: EXIF (Exchangeable Image File Format) and XMP (Extensible Metadata Platform) data provide the history and related details of image editing, while the existence of ICC (International Colour Consortium) or IPTC (The International Press Telecommunications Council) headers serve as sign for an intermediate processing by an Adobe product [14].

This approach consists of the following steps:

1. Read the input image.
2. Extract EXIF data.
3. Check the mismatch between “DateTimeOriginal” and “DateTimeModified” fields.
4. Check the existence of “Adobe Photoshop” in the “Software” field.
5. Extract XMP data.
6. Check the existence of “Adobe Photoshop” in the “CreatorTool” field.
7. Check the existence of ICC header.
8. Check the existence of IPTC header.
9. Alert a manipulation case if any of the 3, 4, 6, 7, and 8 conditions was met, figure (3).

```

<EXIF Metadata>
Adobe header is found.
  Software:  Adobe Photoshop CC 2019 (Windows) (Manipulated Document)

Difference between dates is found, (Manipulated Document)
  Where: DateTimeOriginal:  2008:01:15 19:02:43
        DateTimeModified:  2020:06:24 20:30:40

Make :  HP
Model:  HP psc1500

<XMP Metadata>
Adobe header is found.
  CreatorTool:  Adobe Photoshop CC 2019 (Windows) (Manipulated Document)
Manipulation History: Dates
2020-06-24T19:32:45+03:00
2020-06-24T20:30:40+03:00

ICC header is found : (Manipulated Document)
IPTC header is found : (Manipulated Document)
    
```

Figure (3): output of metadata analysis method. Rounded rectangles indicate the alerts.

2.2 Error Level Analysis (ELA)

Official documents such as certificates and transcripts usually use standard fonts and sizes which make it easy to modify their content. Throughout this section, we modify the ELA method to detect and locate the inserted texts with black gradient colour scheme, as this range of colour is the most prevalence in documents. Our method uses the concept of error level analysis presented by [15]. ELA works by subtracting an image from a compressed version of it at a defined error rate (90%). The image of the difference will highlight the low quality 8 x 8 non-authentic blocks, as these regions degrade at a higher rate at compression [10]. Basically, ELA has high potential in forgery detection of documents, yet there is controversy over its robustness due to the ambiguity of the findings

and reliance on user's interpretation. In our method we build on the image enhancement step in [16] as an input for the detection and localization step. In the last steps we look for the adjacent red, blue or green pixels in the output image, as the presence of these colours distinguish the inserted text regions. Our implementation requires the user's input to select the area containing text.

The steps of the method are:

1. Read the image in RGB format.
2. User selects the text area.
3. Crop the selected part and save it internally in its original and low-level quality.



Figure (4): output of error level analysis method. Images from left to right: original, modified text in green rectangle and the output with marked pixels.

4. Calculate the absolute difference between the corresponding pixel values, for each channel, in the two images from step 3.
5. Enhance the difference image brightness using scaling, and that it is by normalizing the pixel values to be within the “0 – 255” range.
6. Mark the pixels in blue in the resulted image when there is a major difference between the values of the three channels, figure (4).

2.3 Copy-move Forgery Detection

Copy-move tampering is one of the common techniques in the domain of digital images forgery in general [17] and especially scanned documents forgery. This kind of tampering leaves no visual traces, yet it changes the image statistics. By definition, it is the process of copy part of an image (e.g. characters, full words, etc.) and paste it within the same image in order to modify its content. In this method we use the approach which is called “Exact Match” in [3], nevertheless, we propose additional steps to adapt the algorithm to be applied in the field of scanned documents forgery detection. “False Positives”, for example, is one of the top challenging problems encountered by detection algorithms in this application. The almost homogeneous white background of documents and the identical shapes of characters are the key places in which a random duplication detection algorithm may fail and produce false positives. This modification aims to tackle this problem by excluding blocks with prominent whiteness (blocks below a reasonable percentage of blackness, such as: the whitebackground, characters’ outer edges, scanning noise...etc.).

The algorithm consists of the following steps:

1. Read the document image in RGB format.
2. Divide the image into overlapping blocks each with a fixed size of (B * B).
 - Move a sliding window (B * B) over the image from the top-left corner to the bottom-right corner with one-pixel slide for each move.
 - Store the pixel values of each block into a vector.

• The Total Number of Blocks:

$$T = (M - B + 1) * (N - B + 1)$$

where M, N and B are the width, height of the image and the block size respectively.

3. Combine blocks in a single 2D array (feature array) where each row represents a block.
4. Sort the feature array lexicographically by rows, according to their values. As a result, the identical rows will sequence and the matching process later on will be speeded up.
5. Matching: compare between the consecutive rows in the feature array to identify the duplicated blocks.
6. Calculate the maximum whiteness that a block can contain:

$$Max = B^2 * C * 255$$

where B = 4, C = 3, and 255 indicate the block size number of channels, and the maximum pixel value respectively.

7. Calculate the summation of pixel values for each duplicated block that was identified by the Matching step:

$$Sum = \sum_{a=1}^{B^2} \sum_{b=1}^3 Block[a][b]$$

8. Mark blocks that contain whiteness less than 70% of the maximum possible whiteness that a block might have, figure (5):

***ifSum < 0.7 * Max:
thenmarktheblockasduplicated***



Figure (5): output of copy-move detection method. Images from left to right: original, cloned region in green, output with marked repeated block.

3. Experimental Results

As all of the three methods functions independently, we will illustrate the results of each one separately.

3.1 Metadata Analysis Results

Experiments were carried out for the assessment of metadata analysis method. Scanned documents images with “JPG”, “PNG”, and “TIF” file formats were used, where the images with “PNG” and “TIF” formats were converted from “JPG format”. All of the input images were scanned using scanner devices while some of them were edited using Photoshop or Microsoft Paint software. The algorithm’s ability to extract the metadata, and hence detect forged document, varies according type of the image. Table-1 summarizes the algorithm’s performance in extracting each of “EXIF”, “XMP”, “ICC,” and “IPTC” information for each of the three file formats.

	JPG	TIF	PNG
EXIF	Succeeded	Succeeded	Failed
XMP	Succeeded	Succeeded	Succeeded
ICC	Succeeded	Succeeded	Succeeded
IPTC	Succeeded	Succeeded	Failed

Evaluation Metric	Average
Precision	0.745
Recall	0.048
F1 score	0.099

3.2 Error Level Analysis Results

Table-2 summarizes how this approach performs using tampered “JPG” images representing scanned documents. The editing process of those images was performed by “Microsoft Paint” software, where a region was erased and replaced by a new transparent text.

Table-2: performance summary of ELA method.

Where:

$$Precision = \frac{TruePositives}{TruePositives + FalsePositives}$$

$$Recall = \frac{TruePositives}{TruePositives + FalseNegatives}$$

$$F1\ Score = 2 \times \frac{Precision \times Recall}{precision + Recall}$$

True Positives, False Positives and False Negatives are the number of the truly detected blocks, falsely detected blocks and non-detected blocks that must be detected.

3.3 Copy-move Detection Results

In order to assess the modified copy-move forgery detection and localization algorithm in terms of precision, recall, and F1 Score, a Ground-truth must first be defined. For this case, the Ground-truth should indicate the smallest number of overlapping blocks representing the character(s), where each block has at least 30% blackness. The explanation for this is that characters within documents usually have dark colours on a white background, and the “30%” threshold will be sufficient to keep the essential character blocks. Table-3 exemplifies the output of the modified copy-move detection and localization algorithm with different file formats of non-tampered input images. Half of the input images used in the experiment were camera captured while the rest were captured using a scanner. The objective of this experiment is to test the ability of the algorithm to correctly handle the homogeneity exists in the background. The numbers in Table-3 indicate the number of false detected blocks (false positives).

	JPG	TIF	PNG	BMP
Camera Image	48685	47595	47595	47595
Scanner Image	0	0	0	0

Table-3: number of false detected blocks in two camera and scanner sample images.

In another experiment, Table-4 summarizes the performance evaluation of the algorithm for manipulated images of scanned documents. The tampering cases for each image include: identical duplication, duplication with a 20% enlargement, duplication with a 30% enlargement, duplication with a 20% reduction, and duplication with a 30% reduction. The input images were in "TIF" format, where the same experiment were previously carried out using "PNG", "BMP", and "JPG" formats. The results of using the "PNG" and "BMP" file formats were identical to those of using "TIF" format. Contrarily, the algorithm produced surprisingly no true positives and the normal false positives using images with "JPG" format.

Evaluation Metric	Precision	Rec-all	F1 Score
Identical	0.993	1	0.997
Enlargement 20%	0.982	0.26	0.409
Enlargement 30%	0.965	0.116	0.207
Reduction 20%	0.907	0.089	0.161
Reduction 30%	Undefined	0	Undefined

Table-4: performance summary of copy-move detection method for different modifications on the copied region.

4. Discussion

4.1 Metadata Analysis Method

Table-1 clearly shows that converting a document image from "JPG" to "PNG" format can lead to valuable information loss, as "PNG" format does not support EXIF and IPTC as standard tags. On the other side, as "TIF" file format supports EXIF, XMP, ICC, and IPTC tags, the conversion from "JPG" type will not be a problem and all the editing-related information will be preserved. Finally, it is important to note that the efficiency of this approach essentially depends on the assumption that the metadata would be changed by the image editing software.

4.2 Error Level Analysis Method

The need of user's input to select text areas refers to the fact that the modified method fails when processing coloured shapes like stamps and logos. The reason of this is that the coloured pixels of those shapes lead this method to define them as inserted text pixels. Regarding the results of this approach, a good result in terms of precision was obtained. It is important to note that the performance deteriorated significantly when processing some images due to the existence of some coloured pixels in characters. The significant low Recall metric in Table-2 indicates the detection of only a few of the inserted pixels. Despite the low recall, the overall performance still be sufficient when taking the good precision into consideration, since a few correctly detected pixels of a word means that its other pixels are inserted (not-original) too. Another crucial point to mention when using ELA methods in general, is that their performance extremely drops with each re-saving of the image.

4.3 Copy-move Detection Method

Table-3 clearly shows the weakness of the modified algorithm with processing camera captured document images. The high number of false positives is for thousands of identical blocks exist in the background of the document. While the typical background of any scanned document is white, the background of these images was grey due to the poor lighting conditions, figure (6). This fact led to more blackness in the background blocks and thus falsely defining them as characters.

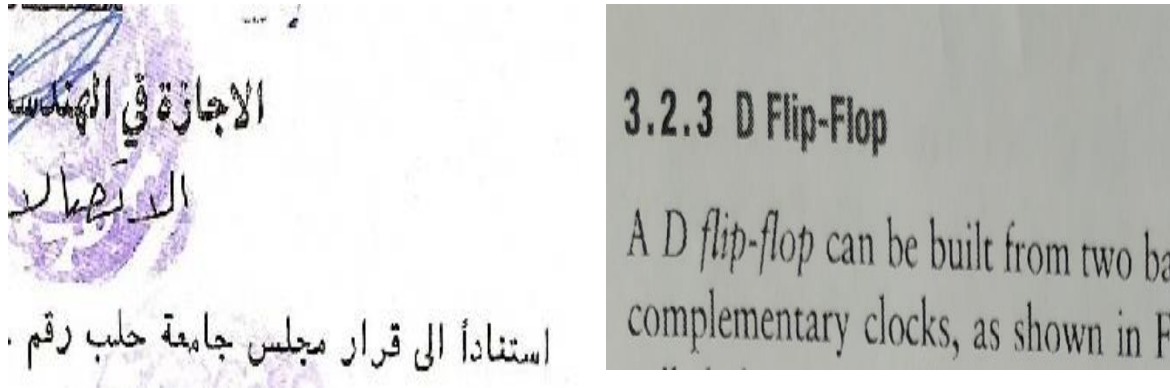


Figure (6): the effect of poor lighting conditions on the input image. From left to right: scanner image, camera image.

As to Table-4, it illustrates how mostly the modified algorithm performs excellently in terms of precision when processing tampered images. This high precision helps in the precise localization of the tampered blocks within the image. The algorithm fails with detecting the duplicated blocks if they were resized to 70%. Regarding of the Recall metric, any applied resizing on the copied region leads to a drastic drop in this metric to below 26%. In other words, less than quarter of the size of a duplicated word is detected. The overall algorithm's performance is still sufficient when considering the high precision, as a few correctly detected blocks of a word implicitly means that the remainder are duplicated too. Another point can also be noticed is that the algorithm is most effective in detecting identical duplications (no resizing).

5. Conclusion

This research presented an end-to-end system for forgery detection in scanned document images using metadata analysis, error level analysis, and copy-move forgery detection methods. Firstly, we worked on the adaption of each the aforementioned methods to the detection of forgeries in scanned document images. After then, and since every method works independently, experiments were performed on each of them to evaluate their performance using specific metrics. Eventually, all of the three methods were combined in a cascade to consist an end-to-end system that detects and localizes the forged regions within the image. The presented results in our research show the possibility to adapt and apply well-known methods to and in the forgery detection in scanned document images. Our future work will focus on detecting overlaid regions by white background, deploying more robust versions of copy-move detection method, and exploiting error level analysis to detect and localize spliced regions. Most importantly, more image processing operations will be performed to modify the test images in order to guarantee the generalization of these three methods.

6. Acknowledgment

This research was carried out as a part of “Documents Authenticity Checker” project by “Syrian Society for Scientific Research – Engineering Clubs”, WWW.SYSSR.ORG.

"For Abroad"

Damascus on 20 January 2020

Signature (Signed)

I, Khudur Sameer Darwish, Mayor of (Almazzeah Jabal Locality), do certify that the details hereinabove stated are true and that Mr. **HUSEN FAISAL ISSA** is the person whose photograph is posted hereinabove, and that I shall notify him with all judicial papers issued for him in case of his absence.

Damascus on 20 January 2020

Mayor Stamp and Signature (Signed and Stamped)

Attestation of the Stamp and Signature (Signed and Stamped)

Attested by the Damascus Governorate on 20 January 2020 (Signed and Stamped)

True translation of the attached document
Damascus on 21/01/2020

Sworn Translator
Muhammad Tarif Alhawari



(1)

"For Bank"

Damascus on 20 January 2022

Signature (Signed)

I, Khudur Sameer Darwish, Mayor of (Almazzeah Jabal Locality), do certify that the details hereinabove stated are true and that Mr. **HUSEN FAISAL ISSA** is the person whose photograph is posted hereinabove, and that I shall notify him with all judicial papers issued for him in case of his absence.

Damascus on 20 January 2022

Mayor Stamp and Signature (Signed and Stamped)

Attestation of the Stamp and Signature (Signed and Stamped)

Attested by the Damascus Governorate on 20 January 2022 (Signed and Stamped)

True translation of the attached document
Damascus on 21/01/2022

Sworn Translator
Muhammad Tarif Alhawari



(2)

Difference between dates is found, (Manipulated Document)

Where: DateTimeOriginal: 2008:01:15 19:02:43
DateTimeModified: 2020:06:24 20:30:40

(3)

"For Bank"

Damascus on 20 January 2022

Signature (Signed)

I, Khudur Sameer Darwish, Mayor of (Almazzeah Jabal Locality), do certify that the details hereinabove stated are true and that Mr. **HUSEN FAISAL ISSA** is the person whose photograph is posted hereinabove, and that I shall notify him with all judicial papers issued for him in case of his absence.

Damascus on 20 January 2022

Mayor Stamp and Signature (Signed and Stamped)

Attestation of the Stamp and Signature (Signed and Stamped)

Attested by the Damascus Governorate on 20 January 2022 (Signed and Stamped)

(4)

"For Bank"

"For Bank"

Damascus on 20 January 2022

Signature (Signed)

I, Khudur Sameer Darwish, Mayor of (Almazzeah Jabal Locality), do certify that the details hereinabove stated are true and that Mr. **HUSEN FAISAL ISSA** is the person whose photograph is posted hereinabove, and that I shall notify him with all judicial papers issued for him in case of his absence.

Damascus on 20 January 2022

Mayor Stamp and Signature (Signed and Stamped)

Attestation of the Stamp and Signature (Signed and Stamped)

Attested by the Damascus Governorate on 20 January 2022 (Signed and Stamped)

(5)

Damascus on 20 January 2022

Signature (Signed)

I, Khudur Sameer Darwish, Mayor of (Almazzeah Jabal Locality), do certify that the details hereinabove stated are true and that Mr. **HUSEN FAISAL ISSA** is the person whose photograph is posted hereinabove, and that I shall notify him with all judicial papers issued for him in case of his absence.

Damascus on 20 January 2022

Mayor Stamp and Signature (Signed and Stamped)

Attestation of the Stamp and Signature (Signed and Stamped)

Attested by the Damascus Governorate on 20 January 2022 (Signed and Stamped)

True translation of the attached document
Damascus on 21/01/2022

Sworn Translator
Muhammad Tarif Alhawari



(6)

Figure (7): scanned document (1), modification (cloning and inserting text) (2), metadata method output (3), ROI selection by user (4), ELA method output (5), copy-move method output (6)

References

1. Lackovic, N. (2020) Thinking with digital images in the post-truth era: a method in critical media literacy.
2. Bayram, S., Sencar, H. T. and Memon, N. (2008) A survey of copy-move forgery detection techniques.
3. Fridrich, J., Soukal D. and Lukas J. (2003) Detection of copy-move forgery in digital images.
4. Mushtaq, S. and Mir, A. H. (2014) Digital image forgeries and passive image authentication techniques: a survey.
5. Abramova, S. and Böhme, R. (2016) Detecting copy-move forgeries in scanned text documents.
6. Gangwar, D. P. and Pathania, A. (2018) Authentication of digital image using Exif metadata and decoding properties.
7. Kee, E., Johnson, M. K. and Farid, H. (2011) Digital image authentication from JPEG headers.
8. Jeronymo, D. C., Borges, Y. C. C. and Coehlo, L. (2017) Image forgery detection by semi-automatic wavelet soft-thresholding with error level analysis.
9. Sudiatmika, I. B. K. and Rahman, F. (2019) Image forgery detection using error level analysis and deep learning.
10. Yancey, R. E., Matloff, N. and Thompson, P. (2019) Multi-stream faster RCNN with ELA for image tampering detection.
11. Carrazoni Entenza, P. (2019) Copy-move forgery detection in scanned text documents using Zernike and Hu moments.
12. Cruz, F., Sidere, N., Coustaty, M., d'Agency, V. P. and Ogier, J (2017) Local binary patterns for document forgery detection.
13. Sencar, H. T. and Memon, N. (2007) Overview of state-of-the-art in digital image forensics
14. Krawetz, N. (2012) Hacker factor.
15. Krawetz, N. (2017) A picture's worth... digital image analysis and forensics.
16. Ciro S. Costa (2015) GitHub.
17. Farid, H. (2009) Image forgery detection, a survey.