# A Two-Step Method for Ensuring Printed Document Integrity using Crossing Number Distances*

Felix Yriarte
*Univ Lyon 2, CNRS, LIRIS,*
F-69676 Bron, France
felix.yriarte@liris.cnrs.fr,

Pauline Puteaux
*Univ. Lille, CNRS, CRIStAL,*
Lille, France
pauline.puteaux@cnrs.fr

Iuliia Tkachenko
*Univ Lyon 2, CNRS, LIRIS,*
F-69676 Bron, France
iuliia.tkachenko@liris.cnrs.fr

*Abstract*—Nowadays, with the use of photo-editing software being mainstream, document integrity verification has become crucial. Indeed, a printed and scanned document undergoes geometric transformations, as well as the addition of black spots, not to mention a decrease in color intensity. The relevant features of an original document, which will be matched against a query document, are stored to be used as a template. We propose a 2-step method that compares a template with a query document to ensure that the query document has not been tampered with. Our method first reverts geometric transformations the document underwent, and then extracts the crossing numbers in that image. A Euclidean distance based matching method is applied to the two sets of crossing numbers, and abnormally distant point groups are flagged as potentially modified. A second step in our method is then applied to analyze the statistical properties of these distance values, to ensure that the document has not been altered. Our results when we apply our method to a database containing administrative documents and tampered versions of these documents – all of which underwent a print and scan process – show the validity of our considerations.

*Index Terms*—document integrity check, print-and-scan process, printed document, document forgery.

## I. INTRODUCTION

Although a number of administrative documents are still distributed in paper form, most of these are then transmitted over the internet, in digital form. Verifying the integrity of digital documents is not a simple task, particularly because of how easy it is to use powerful image editing tools (like Photoshop or GIMP). Additionally, deep learning technologies have recently been used to produce high quality document forgeries [14].

Document hashing works well for digital documents, for example by using OCR (Optical Character Recognition) techniques paired with cryptographic hashing functions [13]. However, OCR stability and accuracy significantly drop if the document is printed and scanned once [1], and plummet after a double print-and-scan (P&S) process [12].

Printed document integrity check can also be done using printer forensics techniques. In this approach, the specific features such as noise intensity, contour roughness, and average gradient of character edges are analyzed as to identify the forged areas [6]. The main drawback of such approaches is the necessity of having knowledge about the printer and scanner used.

Document forgery detection can also be done by constructing document-specific hashes. In this approach, features are extracted from each character, then encoded so that the resulting codes can be used for integrity check [11]. Several of these features – which are also used in biometrics – were shown to be robust to the P&S process [5]. In this work, we deal with administrative document falsification. A genuine document is generated by the authorities. The document's signature is computed and stored in the authority's database. To compute that signature, the coordinates of the features considered in [5] are extracted. That signature can also be stored in a barcode, that could even be integrated to the document. The document could still be used either in digital form, or as a printed and scanned document. When the user transmits the document to any entity, its integrity can be verified by comparing its signature with the stored signature. An overview of the studied document life cycle is illustrated in Fig. 1, both for a genuine and an attack scenario.

The rest of the paper is organized as follows. We present the proposed document signature extraction system in Section II and document integrity check in Section III. The experimental results are presented in Section IV. Finally, we conclude this work and discuss future prospects in Section V.

## II. DOCUMENT SIGNATURE EXTRACTION

We propose a method for document integrity verification that is robust to the P&S process. This method consists of a signature extraction step and of a verification step. The proposed integrity check system is illustrated in Fig. 2. The signature extraction step consists of 1) document image pre-processing and 2) character feature extraction. The integrity check step is applied to identify fields that present abnormal features, and then to confirm or deny that they are indeed forged.

In this section we overview the necessary pre-processing operations and present the features used. The proposed feature matching method is described in Section III.
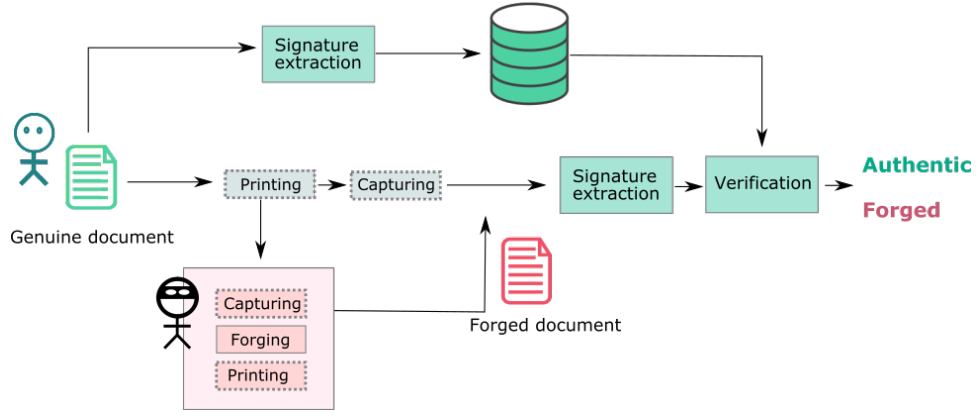
Fig. 1. Overview of the considered scenario: the green blocks correspond to authorities, the red blocks correspond to a forger. The dashed parts represent optional processes.
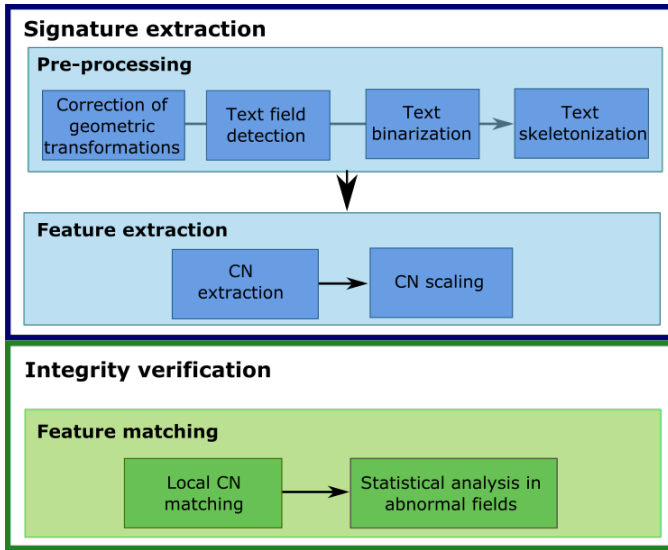


Fig. 2. Pipeline of the document image pre-processing steps and text integrity verification.

## A. Pre-processing operations

The P&S process adds different types of degradations to a document: geometric transformations, addition of black spots, gray-level value change (P&S images are grayscale, whereas digital images usually are binary), as well as compression artifacts that appear after the scan [9]. In order to eliminate part of the deterioration that is due to the P&S process, several pre-processing operations are applied to a document image. Here we list the pre-processing operations that are shown in our document integrity check pipeline illustrated in Fig. 2.

- **Correction of geometric transformations.** Scanning a printed document adds different geometric transformations to the document image, such as rotation, translation, scale change, as well as crop. Even if the scanning operation is done carefully, the resulting image is bound to undergo geometric transformations.
- **Text content detection.** We then apply a text detection step so as to remove spots, and non-textual image contents, such as tables, or the document frame. Indeed, after a P&S process, it is common to find small black spots on the document image. These spots usually are smaller than a character.
  We apply the edge-detection method proposed by Suzuki and Keiichi [10], and filter out unusually small or large content. This allows us to shape an image mask – a binary image of the same resolution as the document image – indicating the fields that are to be considered.
- **Binarization and skeletonization.** In order to extract the feature points from a document image, it is necessary to apply a skeletonization step to it. A skeletonization is a shape thinning operation that returns a 1 pixel-wide skeleton that preserves the connectivity of components. We use the method proposed by Lee in [2]. Before the skeletonization, the image is thresholded using the thresholding method proposed by Otsu [3].

After the application of these pre-processing operations, a skeleton image is obtained, and features can be extracted.

## B. Feature extraction

Since the P&S process induces noise, it is important to take into account features that are robust to such operations. Crossing Numbers (CN) are used in biometrics (namely, fingerprint recognition) as features because of their stability through acquisition noise. They can however also be used to identify textual characters [5], [11]. A skeleton pixel's crossing number value is the number of neighbouring skeleton pixel it has. This gives an indication on the connectivity of every skeleton pixel. We do not consider CN values of 2 for our matching step, since they cannot be used to distinguish different characters.

Furthermore, some characters can present serifs – small lines or strokes attached to the end of longer strokes. Serifs cause the addition of two types of crossing numbers: a bifurcation and up to two ending points. At lower resolutions, serifs are not reliably identified and can be missed ; in which case the matching step fails for CN feature points extracted

from serifs. Since these features are not significant for the matching step, we remove them from the feature point set during the CN feature point extraction. Depending on the scanning quality used, the resulting P&S image's resolution differs. It is thus necessary to scale the extracted CN feature points so as to normalize the data. However, we do not want to induce interpolation errors by applying a scaling operation to the image, which could add discontinuities to a skeleton image. That is why we scale the coordinates of the extracted CN feature points, rather than scaling the document image itself. To do so, every extracted CN feature point's coordinates are multiplied by the ratio $\frac{\text{Reference\_Image\_Resolution}}{\text{Query\_Image\_Resolution}}$. Note that every reference document – since they are numeric – has the same image resolution. We consider that resolution known when scaling the extracted CN feature points.

## III. DOCUMENT INTEGRITY CHECK

After the CN feature points are extracted from the query document image, they are matched against the reference document image feature points (stored in a database as a signature), so as to verify the integrity of the query image. In this section, we introduce the two-step document integrity verification method we propose.

The first step of this method identifies abnormal fields – considered as possibly forged – for which no good candidate features are found in the reference. The second step of our method determines whether these abnormal fields are falsified, or highly distorted only.

### A. Local CN matching

Document falsifications often consist of minor changes (like name, surname, or a salary amount). Therefore, global integrity check approaches do not work well. Fields for which the matching distance is abnormally high are more likely to contain falsifications. However, it is also possible that they are only subject to geometric transformations. These transformations are not always perfectly corrected during the pre-processing step, and can cause slightly increased matching distances for some fields – mostly the document's borders. For example, as shown in Fig. 3.a, the matched sequences are identical, yet the matching distance for each CN feature point is larger than zero. Because of that, some fields can show higher than average matching distances, although they do not contain altered characters. However, we propose a second step where the distribution of matching distances in a field is analyzed, so as to differentiate falsified character sequences and characters that underwent heavy geometric transformations.

### B. Statistical analysis in abnormal fields

We formulate the following hypothesis: if two identical but shifted character sequences are matched against each other, the matching distances are almost similar, as illustrated in Fig. 3.a. The histogram representing the distribution of these matching distances shows low value dispersion (Fig. 3.b). However, when matching against a falsified sequence, the matching
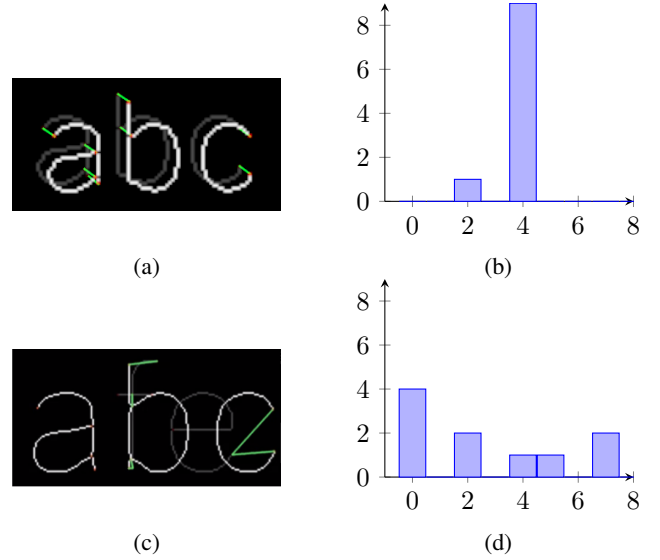


Fig. 3. Example of a) identical shifted sequences of characters (authentic document), b) histogram of distances associated to (a), which illustrates weak dispersion, c) different sequences of characters (forged document), d) histogram of distances associated to (c), which illustrates high dispersion.

distances are relatively random (Fig. 3.c-d), as long as the distributions are analyzed locally – considering the matching distances of the CN feature points inside a field, for instance. If the matching distance distribution within a field is uneven, it is unlikely to contain falsified characters. We can not however say that a distribution that is close to uniform necessarily is falsified, which is why we first locate fields containing high matching distances, and then check whether their distance distribution is uniform or not. A field that both contains high matching distances and of which the distribution is close to uniform is likely to be falsified. If it contains high matching distances of which the distribution is uneven however, it is not likely to contain altered character sequences. To evaluate the distribution of distances within a field, we use the Shannon entropy [7] value of the matching distances $H(D)$, so as to get an idea on whether a distribution is uneven or close to uniform:

$$H(D) = - \sum_{i=0}^{l-1} p(d_i) \log_2(p(d_i)), \qquad (1)$$

where $D$ is the set containing all the distances within a field, of size $k$, and $p(d_i)$ is the probability of occurrence of a distance value $d_i$ ($0 \le d_i < l$).

One can note that the highest distance value $l-1$ depends on the document used. In order to have normalized entropy values (between 0 and 1 bit), we then divide $H(D)$ by the maximal entropy value, which is equal to $\log_2(\min(k, l))$ [4].

## IV. EXPERIMENTAL RESULTS

In this section, we first describe the document database we have considered to perform our experimental results (Section IV-A). Then we present the obtained results by applying our proposed document integrity check method to this database. Finally, we provide a discussion on the limitations

and the drawbacks of our method and give some directions that can be investigated in future work.

### A. Document database used

The Payslip database was proposed by Sidere *et al.* in [8]. It is composed of 200 digitally generated payslips, using names, first names, and addresses among the most common in France. The information presented in these bulletins is therefore fictitious (it is not supposed to represent real people), but close to reality (the different values of the fields are taken from the real administrative documents). The database also contains several falsified versions of every document. The database is made up of 200 genuine documents, and 477 falsified documents. These falsifications were however carried out on digital documents, or on documents that were very slightly altered by a P&S process. We therefore printed and scanned a subset of these documents (both genuine and altered versions) in order to verify the robustness of our method to the P&S process.

We have considered documents with Arial font and a font size of 10 only. The choice of these font type and size are random. We believe that the proposed integrity check works for any font type and size. The subset of image documents used for our experiments[1] is detailed in Table I.

For our experiments, we use 10 different numeric documents as reference documents. We only consider query documents that have gone through a P&S process, as the first step of our proposed method is sufficient to verify the integrity of a numeric document.

|  | Genuine | Forged |
|---|---|---|
| P&S 300 dpi | 10 | 21 |
| P&S 600 dpi | 10 | 21 |
| Double P&S 600 dpi | 10 | 21 |

TABLE I
DESCRIPTION OF THE DATABASE USED FOR OUR EXPERIMENTS.

The implementation of our method was done using Python and standard image processing libraries, such as OpenCV, matplotlib, and scikit-imagenote.

## V. CONCLUSION

In this paper, we propose an efficient method for document integrity verification, that is robust to the print-and-scan process. As a genuine document is created and transmitted, its signature is computed and stored in a database: the document is pre-processed, and its crossing numbers are stored, making up that signature. If a possibly printed-then-scanned document poses as the first one, it is pre-processed, its CN feature points are extracted, and then matched against the stored signature. We consider the Euclidian distance between CN feature points' coordinates, and then their local distribution, so as to identify potential forgeries. Our proposed method is able to both identify a forged document, and locate the field

containing the falsification. Our method is particularly efficient when considering documents that have only been printed-and-scanned once, with a resolution of either 300 or 600 dpi.

In future work, we would like to explore the use of feature augmentation, so as to identify falsifications operated on characters that contain no CN feature points. The identification of the P&S resolution used for a query document would also be an interesting track to investigate, as it would allow for the consideration of different pre-processing methods and thresholds, depending on the P&S resolution used.

### REFERENCES

[1] S. Eskenazi, P. Gomez-Krämer, and J-M. Ogier. A study of the factors influencing ocr stability for hybrid security. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 9, pages 3–8. IEEE, 2017.

[2] C. C. Lee. A Sequential Thinning Algorithm For Image Skeletonization. In Andrew G. Tescher, editor, *Applications of Digital Image Processing VI*, volume 0432, pages 287 – 291. International Society for Optics and Photonics, SPIE, 1984.

[3] N. Otsu. A threshold selection method from gray-level histograms. *IEEE Transactions on Systems, Man, and Cybernetics*, 9(1):62–66, 1979.

[4] P. Puteaux and W. Puech. Noisy encrypted image correction based on shannon entropy measurement in pixel blocks of very small size. In *2018 26th European Signal Processing Conference (EUSIPCO)*, pages 161–165. IEEE, 2018.

[5] P. Puteaux and I. Tkachenko. Crossing number features: from biometrics to printed character matching. In *International Conference on Document Analysis and Recognition*, pages 437–450. Springer, 2021.

[6] S. Shang, N. Memon, and X. Kong. Detecting documents forged by printing and copying. *EURASIP Journal on Advances in Signal Processing*, 2014(1):1–13, 2014.

[7] C. E. Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.

[8] N. Sidere, F. Cruz, M. Coustaty, and J-M. Ogier. A dataset for forgery detection and spotting in document images. In *2017 Seventh International Conference on Emerging Security Technologies (EST)*, pages 26–31. IEEE, 2017.

[9] K. Solanki, U. Madhow, BS. Manjunath, S. Chandrasekaran, and I. El-Khalil. Print and scan resilient data hiding in images. *IEEE Transactions on Information Forensics and Security*, 1(4):464–478, 2006.

[10] S. Suzuki and Keiichi A. Topological structural analysis of digitized binary images by border following. *Computer vision, graphics, and image processing*, 30(1):32–46, 1985.

[11] L. Tan and X. Sun. Robust text hashing for content-based document authentication. *Information Technology Journal*, 10(8):1608–1613, 2011.

[12] I. Tkachenko and P. Gomez-Krämer. Robustness of character recognition techniques to double print-and-scan process. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 9, pages 27–32. IEEE, 2017.

[13] R. Villán, S. Voloshynovskiy, O. Koval, F. Deguillaume, and T. Pun. Tamper-proofing of electronic and printed text documents via robust hashing and data-hiding. In *Security, Steganography, and Watermarking of Multimedia Contents IX*, volume 6505, pages 633–644. SPIE, 2007.

[14] L. Zhao, C. Chen, and J. Huang. Deep learning-based forgery attack on document images. *IEEE Transactions on Image Processing*, 30:7964–7979, 2021.

[1]To stimulate a collaboration and reproducible results, the python code as well as the augmented Payslip database are publicly available via this link: https://gitlab.liris.cnrs.fr/gdr_isis_fuzzydoc/fuzzydoc_cn_distances