

# OCR Graph Features for Manipulation Detection in Documents

Hailey Joren, Otkrist Gupta, Dan Raviv

hjoren@ucsd.edu, otkrist.gupta@lendbuzz.com, dan.raviv@lendbuzz.com

Lendbuzz

125 High Street

Boston, Massachusetts 02110

## Abstract

Detecting manipulations in digital documents is becoming increasingly important for information verification purposes. Due to the proliferation of image editing software, altering key information in documents has become widely accessible. Nearly all approaches in this domain rely on a procedural approach, using carefully generated features and a hand-tuned scoring system, rather than a data-driven and generalizable approach. We frame this issue as a graph comparison problem using the character bounding boxes, and propose a model that leverages graph features using OCR (Optical Character Recognition). Our model relies on a data-driven approach to detect alterations by training a random forest classifier on the graph-based OCR features. We evaluate our algorithm's forgery detection performance on dataset constructed from real business documents with slight forgery imperfections. Our proposed model dramatically outperforms the most closely-related document manipulation detection model on this task.

## Introduction

The use of digital documents for information verification is becoming increasingly common in domains such as finance, insurance, and administration. For example, digital documents may be used for the purpose of filing insurance claims, verifying address or income, or expense management. As these use cases increase, opportunities for deception through document manipulation naturally follow, as successful forgery of these types of documents can lead to personal reward. Recent technological advances have made manipulation of digital documents easier than ever, allowing forgers to cheaply manipulate key portions of a document. These portions may include information such as account balance, address, name, or salary. Using simple image editing software, forgers can alter a few characters or words in the document to achieve their purposes. This type of forgery is particularly difficult to mitigate as these types of documents don't contain extrinsic artifacts such as watermarks for authentication purposes. Any method for detecting manipulation in these types of documents must thus rely on the intrinsic properties of the document.

Digital document manipulation occupies a unique place within image manipulation detection, as properties such as color or texture that are normally leveraged for manipulation detection may not be as useful in digital documents. For example, many digital documents contain little or no color, and often an entirely white background. In addition, there are a variety of mechanisms available to alter a character or word in the document. The forger can copy and paste from another portion of the document, splice from a separate document, or construct the character from scratch by changing individual pixels (see Figure 1).

Manipulation detection methods for both images and documents rely on imperfections in the forgery process. In documents, these imperfections are often found in the alignment and sizing of the forged character (van Beusekom and Shafait 2011). We propose a method for detecting these forgery imperfections in document manipulation using OCR (Optical Character Recognition) graph features. We frame the problem of manipulation detection as a graph comparison problem, in which each character in the document is considered as the central node of a sub-graph within the larger graph of the document. Each sub-graph contains information about the central node and its neighbors, including features such as OCR box size, distance from the central node and alignment. We then train a random forest to recognize sub-graphs in which the central node has been manipulated.

We construct a dataset which directly measures the ability of our model to detect these types of imperfections by stochastically shifting or scaling characters documents. Constructing this type of dataset constrains the model to rely on forgery imperfections, as texture or pixel-based features may be distorted or removed through post-processing like printing and scanning, Gaussian blurring or JPEG compression (Bayar and Stamm 2016).

Our contributions can be summarized as follows:

- We are the first to consider the problem of document manipulation as a graph comparison problem and to leverage OCR features as graph features for manipulation detection. Our proposed method significantly outperforms the most closely related model designed for document manipulation detection.
- We present a data-driven approach that uses a random

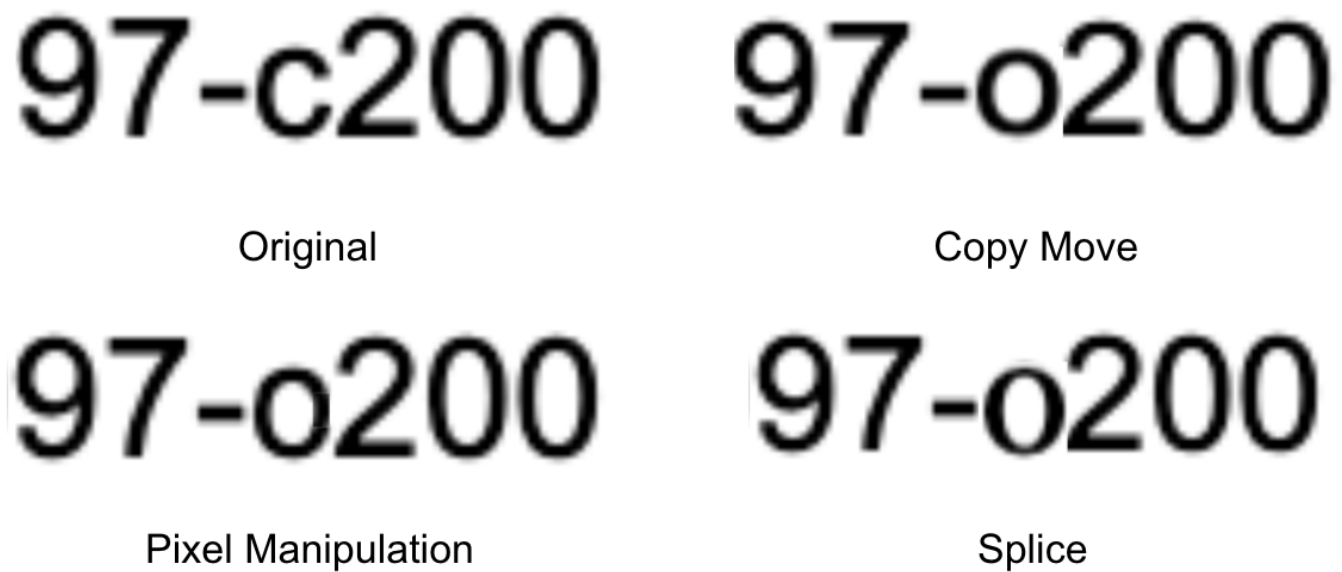


Figure 1: Three possible manipulation types in document forgery include pixel manipulation, copy-move, and splicing. In pixel manipulation, individual pixels are manipulated to change the character. In this case, the gap in the 'c' is filled in to become an 'o'. In copy-move, the character is copied and pasted from another part of the document. In splicing, the character is copied and pasted from another document. In this case, font differences can be observed under careful inspection.

forest architecture, in contrast with the procedural approaches previously employed for detecting forgeries.

- We evaluate our algorithm on a dataset that directly measures the ability of a manipulation detection model to detect forgery imperfections.

### Related Work

#### Manipulation Detection in Images

Digital document manipulation detection occupies a space within image manipulation detection more broadly and builds upon work in this area. For example, (Fridrich, Soukal, and Lukas 2003) investigate detection of copy-move forgery in images using block matching methods. (Lin et al. 2009) develop a method for detecting copy-move forgery in images using Hu moment features on blocks of pixels. (Popescu and Farid 2005) use a pixel-based method for re-sampling using the statistical artifacts created during the manipulation operation. (Stamm and Liu 2010) use a similar approach for detecting contrast enhancement, histogram equalization and JPEG-compression.

However, documents present unique challenges and opportunities when compared with manipulation detection in images. For example, as noted by (Lin et al. 2009), one challenge in detecting copy-move forgeries in images is trying to compare every possible pair of pixel blocks. In contrast, in documents this number of comparisons can be drastically reduced by ignoring white space and only comparing boxes containing the same alphanumeric character. Additionally, methods designed for forgery detection in images that leverage texture or pixel-based statistical properties may perform worse on documents due to large amounts of white space

and less color information, as well as the potential for printing and scanning documents and removing these features. Our proposed model in particular attempts to leverage the structured and predictable nature of documents to develop more discriminative features.

#### Manipulation Detection in Documents

Several methods for detecting forgeries rely on identifying the type of printer or exact machine from which a document was printed (van Beusekom, Shafait, and Breuel 2013a; Shang, Memon, and Kong 2014). In contrast, (Ahmed and Shafait 2014) use intrinsic document elements, or portions of the document that remain the same across different copies of the document, to detect manipulations in documents from the same source. (Bertrand et al. 2015) employ a conditional random field model and font features to detect manipulations in which a character or word has been spliced from a document with different fonts. (Abramova et al. 2016) investigate methods for detecting copy-move forgery in scanned text documents using block-based methods on characters that have been copied onto a "forgery" line at the bottom of the page. (van Beusekom, Shafait, and Breuel 2013b) develop a method using text-line alignment and skew features to determine the authenticity of each line on the page, in contrast with other works (including ours) that examine alterations on a word or character level. (Shang, Kong, and You 2015) use geometric properties of characters to distinguish between documents printed from laser printers, ink printers, and electrostatic copiers for purposes of document authentication. (Cruz et al. 2017) use Otsu binarization to extract patches to compare using local binary pattern features combined with a Support Vector Machine

(SVM) for forgery classification. (Van Beusekom, Stahl, and Shafait 2012) present findings from deploying manipulation techniques on real-world datasets to identify weaknesses. (Chernyshova et al. 2019) propose a system for optical font recognition for the purpose of forgery detection. (Sidere et al. 2017) publish a dataset for evaluating performance on forgery detection in French payslips. (Gupta and Kumar 2020) explore the benefits of using boosting and bagging in forgery detection.

## Methodology

We frame the manipulation detection problem in documents as a graph comparison problem, where each character is the center of a graph containing  $2n + 1$  nodes, consisting of the central character node and  $n$  nodes on each side horizontally. We experiment with various values of  $n$ , ranging from 3-9 (see Table 1). We attempt to include only nodes that are on the same line of text by measuring the difference in the y-values of each node in comparison with the heights. Experimentally, we found that when the difference between the y-values of the top of the OCR boxes was less than .85 of the height of the taller box, the boxes were on the same line. In other words, nodes  $n1$  and  $n2$  are considered to be on the same line when

$$|y_{0n1} - y_{0n2}| < (.85 * \max(\text{height}_{n1}, \text{height}_{n2})) \quad (1)$$

where  $y_{0n}$  indicates the y-value of node  $n$  and  $\text{height}_n$  corresponds to the height of node  $n$ . When the central node fell near the left or right end of a line, we imputed the missing values with the values from the other side of the central node when possible.

Each of the features of each node is concatenated into a single vector to describe the sub-graph for each character in the document. For training, we assign the label 1 when the central node has been manipulated, and 0 when the central node is pristine or unaltered. We then train a random forest on the dataset.

The pseudo-code for reproducing our model is provided in the supplementary materials.

## Features

For each node in the sub-graph, we construct a set of features to describe the node, including:

- **Height and width:** Tesseract OCR provides the bounding box of each character, from which the approximate height and width of the character can be determined. If a character has been manipulated, imperfections in the manipulation often manifest in a slight character size difference.
- **Y-value Difference:** We take the relative displacement of neighbouring nodes along the y-axis. We consider the uppermost y-value of the central node as 0, and compute the positive or negative difference of the neighboring node y-values with respect to the central node. This feature is useful as a proxy for alignment error, as exact alignment of the forged character is difficult to achieve even with digital image manipulation tools.

- **Distance:** The Euclidean distance from the center of the each node to the center of the central node is calculated. Information about distance from the central node allows greater weighting of nodes that are closer to the central node. The distance is formulated as

$$d((x_1, y_1), (x_2, y_2)) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \quad (2)$$

where  $(x_1, y_1), (x_2, y_2)$  refer to the center of the central node character box and the center of a neighboring node character box respectively.

- **Hu Moments:** The seven Hu moments are invariate to rotation, scale and translation and allow for information about the character value within the box (Hu 1962; Huang and Leng 2010). The seven Hu moments are defined as

$$\begin{aligned} \phi_1 &= \eta_{20} + \eta_{02} \\ \phi_2 &= (\eta_{20} - \eta_{02})^2 + 4\eta_{11}^2 \\ \phi_3 &= (\eta_{30} - 3\eta_{12})^2 + (3\eta_{21} - \mu_{03})^2 \\ \phi_4 &= (\eta_{30} - \eta_{12})^2 + (\eta_{21} - \mu_{03})^2 \\ \phi_5 &= (\eta_{30} - 3\eta_{12})(\eta_{30} + \eta_{12})[(\eta_{30} + \eta_{12})^2 - 3(\eta_{21} + \eta_{03})^2] \\ &\quad + (3\eta_{21} - \eta_{03})(\eta_{21} + \eta_{03})[3(\eta_{30} + \eta_{12})^2 - (\eta_{21} + \eta_{03})^2] \\ \phi_6 &= (\eta_{20} - \eta_{02})[(\eta_{30} + \eta_{12})^2 - (\eta_{21} + \eta_{03})^2] \\ &\quad + 4\eta_{11}(\eta_{30} + \eta_{12})(\eta_{21} + \eta_{03}) \\ \phi_7 &= (3\eta_{21} - \eta_{03})(\eta_{30} + \eta_{12})[(\eta_{30} + \eta_{12})^2 - 3(\eta_{21} + \eta_{03})^2] \\ &\quad - (\eta_{30} - 3\eta_{12})(\eta_{21} + \eta_{03})[3(\eta_{30} + \eta_{12})^2 - (\eta_{21} + \eta_{03})^2] \end{aligned} \quad (3)$$

where each normalized centroid moment  $\eta_{pq}$  is defined as

$$\eta_{pq} = \frac{\mu_{pq}}{\mu_{00}^\gamma}, \gamma = (p + q + 2)/2, p + q = 2, 3, \dots \quad (4)$$

and each centroid moment  $\mu_{pq}$  is computed as

$$\mu_{p,q} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \bar{x})^p (y - \bar{y})^q f(x, y) dx, dy \quad p, q = 0, 1, 2, \dots \quad (5)$$

where the pixel point  $(\bar{x}, \bar{y})$  is the centroid of the image, or in our case, the centroid of the character bounding box.

- **Principal Inertia Axis:** The Principal Inertia Axis is obtained through Singular Value Decomposition (SVD) on the  $\mu_{20}$ ,  $\mu_{11}$  and  $\mu_{02}$  Hu moments.

## Experiments

As described above, we construct a random forest using Python (Van Rossum and Drake Jr 1995) and the sci-kit learn library (Pedregosa et al. 2011). We ran randomized hyper-parameter search for 480 iterations on search space of the parameters listed in Table 1. Each experiment was run on a cross-validation set of five folds, and the average performance and standard deviation of the folds are reported in the Results section. The experiments were performed on an Intel Core i7-8700 machine with x86-64 architecture and 12 cores.

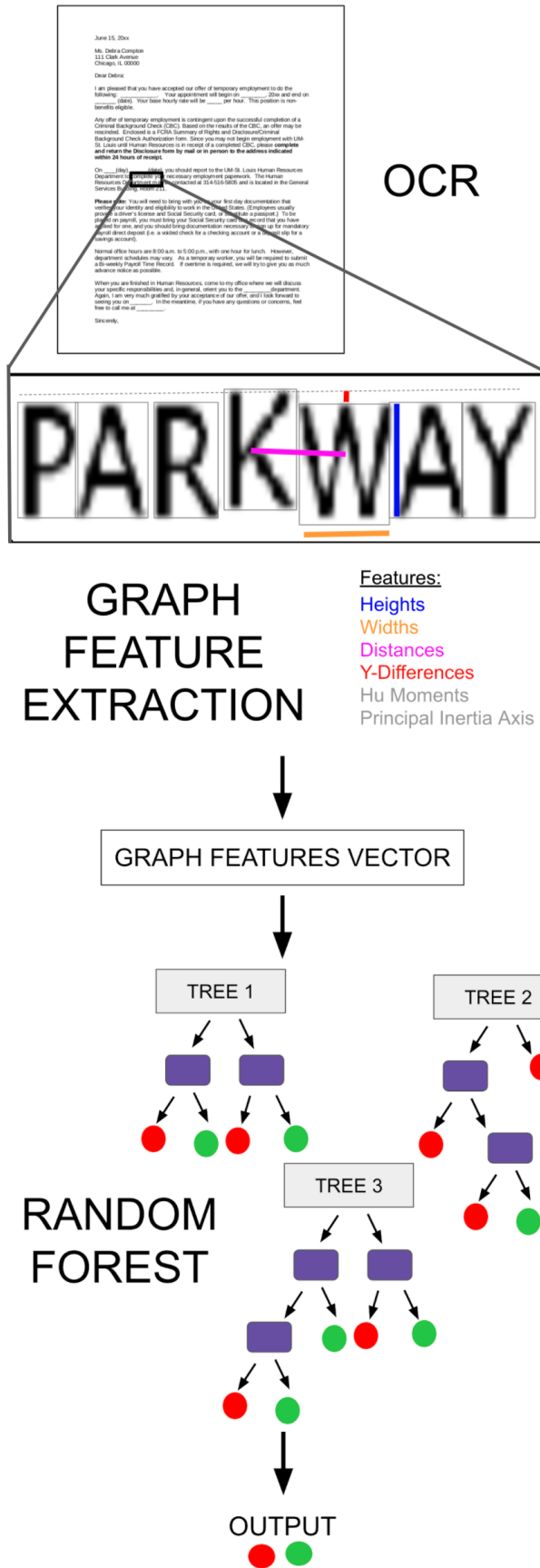


Figure 2: Our pipeline begins with extracting character bounding boxes using Optical Character Recognition (OCR) via Tesseract. The boxes are then used to construct graph features for each potentially manipulated character. In this figure, the 'K' is the character being examined for manipulation. Hu moments and principal inertia axis features are not shown. The graph features are then inputted into a trained random forest, which classifies the character as manipulated or pristine.



Figure 3: The sub-graph contains the central node (red, 'P'), or the OCR character box being inspected for manipulation artifacts, and  $n$  neighbors on either side (blue). Features include box characteristics such as height, width, respective y-value, distance from the central node, Hu moments, and principal inertia axis.

Number of trees	1000, 250, 1500, 1750, 2000, 2250, 2500, 2750, 3000, 3250, 3500, 3750, 4000, 5000
Max tree depth	5, 10, 15, 20, 25, 30, 35, 40, 45, 50
Minimum number of samples per leaf	4, 6, 8, 10, 12
Number of Neighbors ( $n$ )	3, 5, 7, 9

Table 1: Hyper-parameter ranges for the proposed model

## Dataset

Our dataset consists of 359 finance-related documents, including bank statements, offer letters, credit card statements, bills, and tax returns. We split the documents into train and test sets, with 287 and 72 documents in each set respectively. Each document has one or more pages, together totalling 1470 pages for training and 389 pages for testing.

Each document original is a PDF, for which the PDF text boxes are stored in the document. The PDF is first converted to an image, and as we iterate through the character boxes in the document, we alter the box on the image through either shifting or scaling with probability 0.05. Scaling or shifting is applied stochastically within four ranges of values: shifting 1-5px, shifting 5-10px, scaling 7%-14%, and scaling 15%-25%. Empirically, we found that shifting or scaling by values smaller than these did not change the size or position of the character. Values larger than these produces manipulations that were easily visible to the human eye and often significantly disrupted the OCR extraction.

If the box is altered, we store the coordinates of the altered box in the image for later use as ground truth. The resulting dataset is a collection of PNG images of documents of which 5% of the characters have been altered through scaling or shifting based on the PDF text boxes. Notably, the PDF text boxes are not provided to the manipulation detection model, and are distinct from the character boxes obtained through Tesseract OCR. See Figure 4 for examples of scaling and shifting.

While the original documents used for experimentation cannot be shared due to their sensitive nature, we provide pseudo-code in the supplementary materials for processing the dataset which can be performed on any set of similar real-world business documents.

## Comparison

We compare our model’s performance with the system proposed in (Bertrand et al. 2013), which similarly uses intrinsic

document features on the character level to detect character-level manipulations. The proposed method also uses Tesseract OCR (Smith 2007) for character extraction. The authors combine the results of two separate methods for classification of characters as manipulated or pristine.

In the first method, the authors compute a distance between characters of the same alphanumeric type by constructing feature vectors consisting of the seven Hu moments (Hu 1962; Terrades, Tabbone, and Valveny 2007) and calculating the Euclidean distance. Character pairs are flagged when the distance is either suspiciously small, indicating a copy-move forgery, or suspiciously large, indicating a character that is of a different size or font than the other characters of the same type in the document.

In the second method, the authors propose finding a set of forgery clues grouped by a vector  $W$  for each character. The feature vector  $W$  consists of the character size based on the bounding box as obtained through Tesseract OCR, the character principal inertia axis as obtained through Singular Value Decomposition (SVD) on the  $\mu_{20}$ ,  $\mu_{11}$  and  $\mu_{02}$  Hu moments, and the character horizontal alignment on the line level. The feature vector  $W$  of each character is then compared to a data model  $M_c$  for that character alphabetic class using the Mahalanobis distance, which is generated with the document training set. Similar to the first method, a scoring system determines if the distance is small enough or large enough to arouse suspicion of manipulation.

Finally, the two indicators from the two methods are combined to give each character a score, for a final classification of manipulated or pristine.

This method is an appropriate comparison for our method due to its similar character-based features and reliance on Tesseract OCR for determining character bounding boxes. Specifically, the feature extraction methods used by our model are a subset of those used by the model proposed in (Bertrand et al. 2013). Comparison with this model allows us to disentangle the utility of the feature extraction employed from the combined utility of the graph-based features and the random forest architecture.

We performed random search on a range of values and tune the hyper-parameters for a total of 480 iterations. The ranges of each value are shown in Table 2.

## Results

The results in Tables 3, 4 5, 6 indicate that while the model proposed in (Bertrand et al. 2013) achieves a reasonable F1 score and high recall, it reports a high number of false positives. In contrast, our model is able to achieve impressive

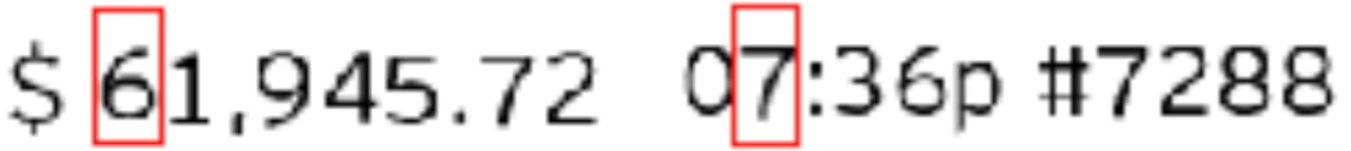


Figure 4: An example of shifting (left) and scaling (right) manipulations. On the left, the six has been shifted up slightly to mimic the y-value alignment artifact introduced during a copy-move, pixel adjusting, or splicing operation. On the right, the seven has been slightly enlarged, as could be similarly expected during manipulation operations.

Method 1 Upper Threshold	80, 95, 90
Method 1 Lower Threshold	0, 5, 10
Mahalanobis Threshold	5, 6, 7
Log Transform Hu Moment Values	True, False

Table 2: Hyper-parameter ranges used for training the model described in (Bertrand et al. 2013)

	Proposed Method	STD	Bertrand 2013	STD
Precision	<b>0.8651</b>	$\pm 0.0134$	0.2927	$\pm 0.0144$
Recall	0.7461	$\pm 0.0181$	<b>0.7862</b>	$\pm 0.0170$
Accuracy	<b>0.9857</b>	$\pm 0.0112$	0.3790	$\pm 0.0128$
F1 Score	<b>0.8012</b>	$\pm 0.0181$	0.4265	$\pm 0.0170$

Table 3: Mean metrics and standard deviations for scaling 5% of characters by a range of 15% to 25% for each of the five splits

performance given even minute manipulations such as shifting by a single pixel or scaling by as little as 7%. We expected a higher F1 score with the "easier" or larger manipulations than was achieved (see Tables 3 and 5), and hypothesize that larger manipulations may disrupt our models ability to correctly distinguish between characters on the same line and characters on different lines.

## Conclusion

We present a system for detecting forgery imperfections using OCR graph features with a random forest architecture. We evaluate this algorithm on detecting small shifting and scaling manipulations in real financial documents. As shown, our model outperforms the most closely related model significantly on this task. It seems that at least for

	Proposed Method	STD	Bertrand 2013	STD
Precision	<b>0.8709</b>	$\pm 0.0143$	0.2495	$\pm 0.0153$
Recall	0.6177	$\pm 0.0221$	<b>0.7971</b>	$\pm 0.0315$
Accuracy	<b>0.9817</b>	$\pm 0.0212$	0.3506	$\pm 0.0122$
F1 Score	<b>0.7228</b>	$\pm 0.0210$	0.3800	$\pm 0.0206$

Table 4: Mean metrics and standard deviations for scaling 5% of characters by a range of 7% to 14% for each of the five splits

	Proposed Method	STD	Bertrand 2013	STD
Precision	<b>0.8187</b>	$\pm 0.0155$	0.2111	$\pm 0.0166$
Recall	0.7440	$\pm 0.0162$	<b>0.8047</b>	$\pm 0.0161$
Accuracy	<b>0.9854</b>	$\pm 0.0172$	0.3318	$\pm 0.0146$
F1 Score	<b>0.7796</b>	$\pm 0.0102$	0.3343	$\pm 0.0106$

Table 5: Mean metrics and standard deviations for shifting 5% of characters by a range of 5-10 pixels for each of the five splits

	Proposed Method	STD	Bertrand 2013	STD
Precision	<b>0.8604</b>	$\pm 0.0101$	0.2498	$\pm 0.0091$
Recall	0.6666	$\pm 0.012$	<b>0.7896</b>	$\pm 0.014$
Accuracy	<b>0.9829</b>	$\pm 0.0198$	0.3494	$\pm 0.0111$
F1 Score	<b>0.7512</b>	$\pm 0.0121$	0.3795	$\pm 0.0116$

Table 6: Mean metrics and standard deviations for shifting 5% of characters by a range of 1-5 pixels for each of the five splits

documents, a data-driven approach may be required instead of a procedural approach for obtaining satisfactory manipulation detection results. In addition, it appears that leveraging the structured nature of digital documents by incorporating graph features offers additional performance gains. The unique nature of digital documents, and their importance in manipulation detection research, warrants special attention and inquiry into developing methods for identifying malicious manipulations.

Future work could include leveraging graph anomaly detection methods for a more efficient graph-based detection system. Additionally, the OCR graph features could be combined with pixel-based features from traditional image manipulation detection for potential improved performance. We also observed that the character boxes resulting from Tesseract OCR were sometimes inaccurate, failing to accurately enclose each character. Improving the underlying OCR method could further improve the manipulation detection results.

## Ethics Statement

Fraud detection is an active area of research and can be of great benefit to society. Detection and prevention of document fraud can potentially help reduce misinformation and

improve resource allocation. Users and organizations often spend a large amount of resources trying to mitigate fraud, usually through human auditors. Automating fraud detection can help reduce this burden. However, as with any machine learning model, caution should be exercised, particularly when considering acceptable false positive rates. Additional human verification may be required when attempting to detect fraud in real life scenarios. Data driven approaches like ours can help mitigate such risks but may not eliminate them, and further research in this domain should be pursued.

## References

- Abramova, S.; et al. 2016. Detecting copy–move forgeries in scanned text documents. *Electronic Imaging* 2016(8): 1–9.
- Ahmed, A. G. H.; and Shafait, F. 2014. Forgery detection based on intrinsic document contents. In *2014 11th IAPR International Workshop on Document Analysis Systems*, 252–256. IEEE.
- Bayar, B.; and Stamm, M. C. 2016. A deep learning approach to universal image manipulation detection using a new convolutional layer. In *Proceedings of the 4th ACM Workshop on Information Hiding and Multimedia Security*, 5–10.
- Bertrand, R.; Gomez-Kramer, P.; Terrades, O. R.; Franco, P.; and Ogier, J.-M. 2013. A system based on intrinsic features for fraudulent document detection. In *2013 12th International conference on document analysis and recognition*, 106–110. IEEE.
- Bertrand, R.; Terrades, O. R.; Gomez-Kramer, P.; Franco, P.; and Ogier, J.-M. 2015. A conditional random field model for font forgery detection. In *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, 576–580. IEEE.
- Chernyshova, Y. S.; Aliev, M. A.; Gushchanskaia, E. S.; and Sheshkus, A. V. 2019. Optical font recognition in smartphone-captured images and its applicability for ID forgery detection. In *Eleventh International Conference on Machine Vision (ICMV 2018)*, volume 11041, 110411J. International Society for Optics and Photonics.
- Cruz, F.; Sidere, N.; Coustaty, M.; D’Andecy, V. P.; and Ogier, J.-M. 2017. Local binary patterns for document forgery detection. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 1, 1223–1228. IEEE.
- Fridrich, A. J.; Soukal, B. D.; and Lukas, A. J. 2003. Detection of copy-move forgery in digital images. In *in Proceedings of Digital Forensic Research Workshop*. Citeseer.
- Gupta, S.; and Kumar, M. 2020. Forensic document examination system using boosting and bagging methodologies. *Soft Computing* 24(7): 5409–5426. ISSN 1433-7479. doi:10.1007/s00500-019-04297-5. URL <https://doi.org/10.1007/s00500-019-04297-5>.
- Hu, M.-K. 1962. Visual pattern recognition by moment invariants. *IRE transactions on information theory* 8(2): 179–187.
- Huang, Z.; and Leng, J. 2010. Analysis of Hu’s moment invariants on image scaling and rotation. In *2010 2nd International Conference on Computer Engineering and Technology*, volume 7, V7–476. IEEE.
- Lin, H.-J.; Wang, C.-W.; Kao, Y.-T.; et al. 2009. Fast copy-move forgery detection. *WSEAS Transactions on Signal Processing* 5(5): 188–197.
- Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; Vanderplas, J.; Passos, A.; Cournapeau, D.; Brucher, M.; Perrot, M.; and Duchesnay, E. 2011. Scikitlearn: Machine Learning in Python. *Journal of Machine Learning Research* 12: 2825–2830.
- Popescu, A. C.; and Farid, H. 2005. Exposing digital forgeries by detecting traces of resampling. *IEEE Transactions on signal processing* 53(2): 758–767.
- Shang, S.; Kong, X.; and You, X. 2015. Document forgery detection using distortion mutation of geometric parameters in characters. *J. Electronic Imaging* doi:10.1117/1.JEI.24.2.023008.
- Shang, S.; Memon, N.; and Kong, X. 2014. Detecting documents forged by printing and copying. *EURASIP Journal on Advances in Signal Processing* 2014(1): 140.
- Sidere, N.; Cruz, F.; Coustaty, M.; and Ogier, J.-M. 2017. A dataset for forgery detection and spotting in document images. In *2017 Seventh International Conference on Emerging Security Technologies (EST)*, 26–31. IEEE.
- Smith, R. 2007. An overview of the Tesseract OCR engine. In *Ninth international conference on document analysis and recognition (ICDAR 2007)*, volume 2, 629–633. IEEE.
- Stamm, M. C.; and Liu, K. R. 2010. Forensic detection of image manipulation using statistical intrinsic fingerprints. *IEEE Transactions on Information Forensics and Security* 5(3): 492–506.
- Terrades, O. R.; Tabbone, S.; and Valveny, E. 2007. A review of shape descriptors for document analysis. In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, volume 1, 227–231. IEEE.
- van Beusekom, J.; and Shafait, F. 2011. Distortion measurement for automatic document verification. In *2011 International Conference on Document Analysis and Recognition*, 289–293. IEEE.
- van Beusekom, J.; Shafait, F.; and Breuel, T. M. 2013a. Automatic authentication of color laser print-outs using machine identification codes. *Pattern Analysis and Applications* 16(4): 663–678.
- van Beusekom, J.; Shafait, F.; and Breuel, T. M. 2013b. Text-line examination for document forgery detection. *International Journal on Document Analysis and Recognition (IJDA)* 16(2): 189–207. ISSN 1433-2825. doi:10.1007/s10032-011-0181-5. URL <https://doi.org/10.1007/s10032-011-0181-5>.
- Van Beusekom, J.; Stahl, A.; and Shafait, F. 2012. Lessons learned from automatic forgery detection in over 100,000 invoices. In *Computational Forensics*, 130–142. Springer.

Van Rossum, G.; and Drake Jr, F. L. 1995. *Python reference manual*. Centrum voor Wiskunde en Informatica Amsterdam.